



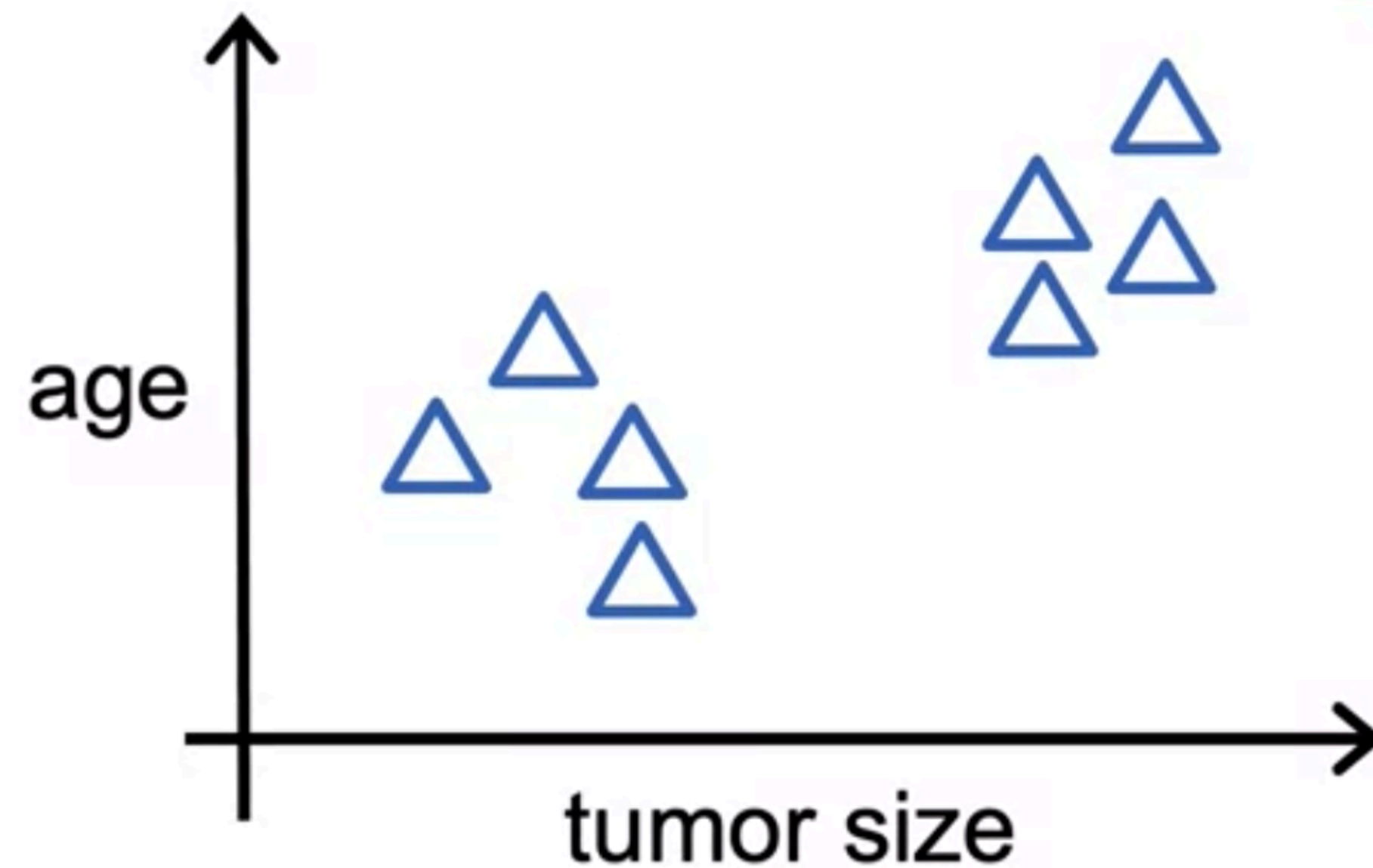
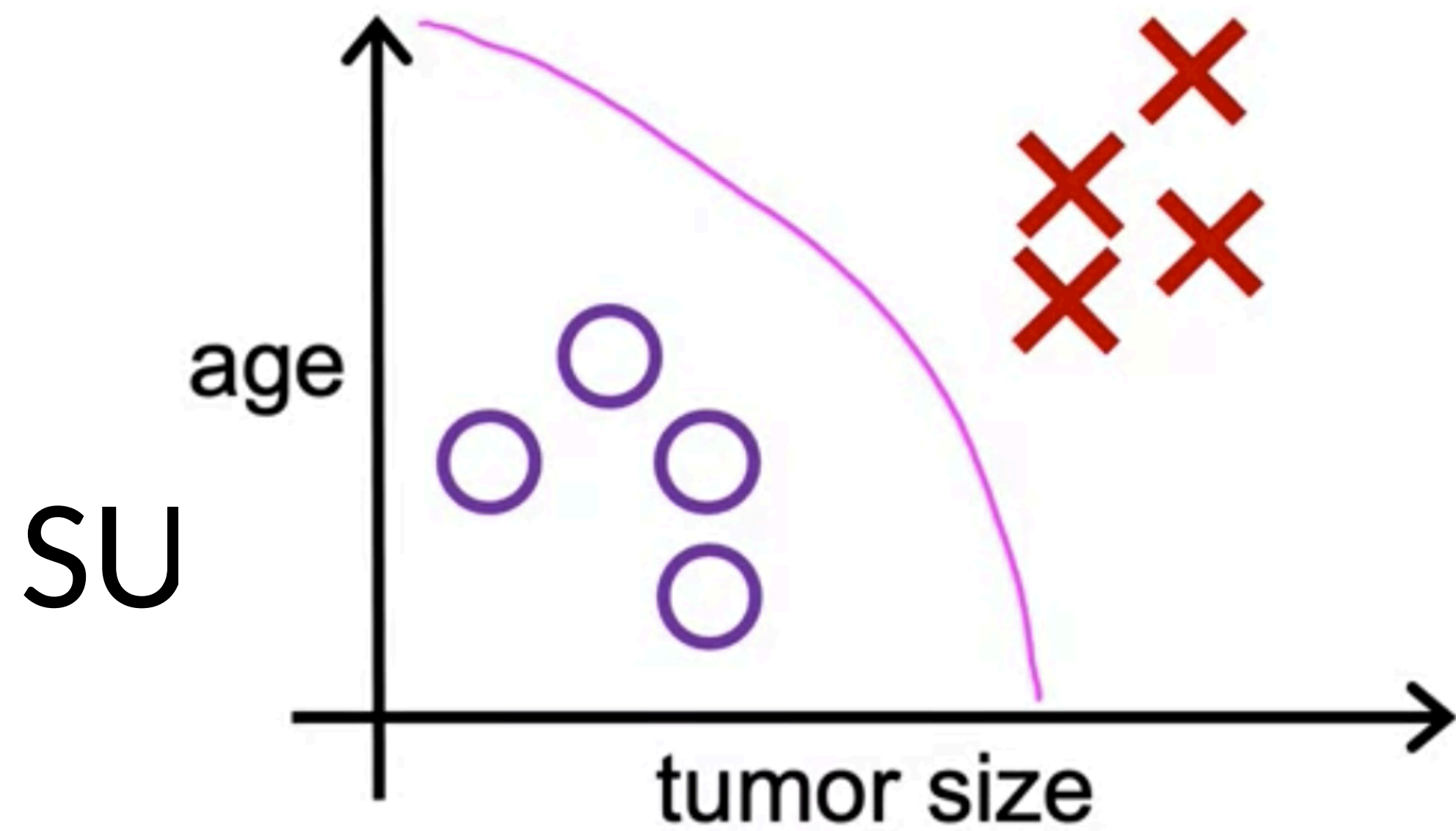
AI Bridge

Lecture 8

UNSUPERVISED LEARNING

Supervised learning
Learn from data labeled
with the "right answers"

Unsupervised learning



TWO MAIN APPLICATIONS

Clustering

Dimension reduction

Clustering: Google news



Giant panda gives birth to rare twin cubs at Japan's oldest zoo

USA TODAY · 6 hours ago

- Giant panda gives birth to twin cubs at Japan's oldest zoo

CBS News · 7 hours ago

- Giant panda gives birth to twin cubs at Tokyo's Ueno Zoo


WHBL News · 16 hours ago

- A Joyful Surprise at Japan's Oldest Zoo: The Birth of Twin Pandas

The New York Times · 1 hour ago

- Twin Panda Cubs Born at Tokyo's Ueno Zoo

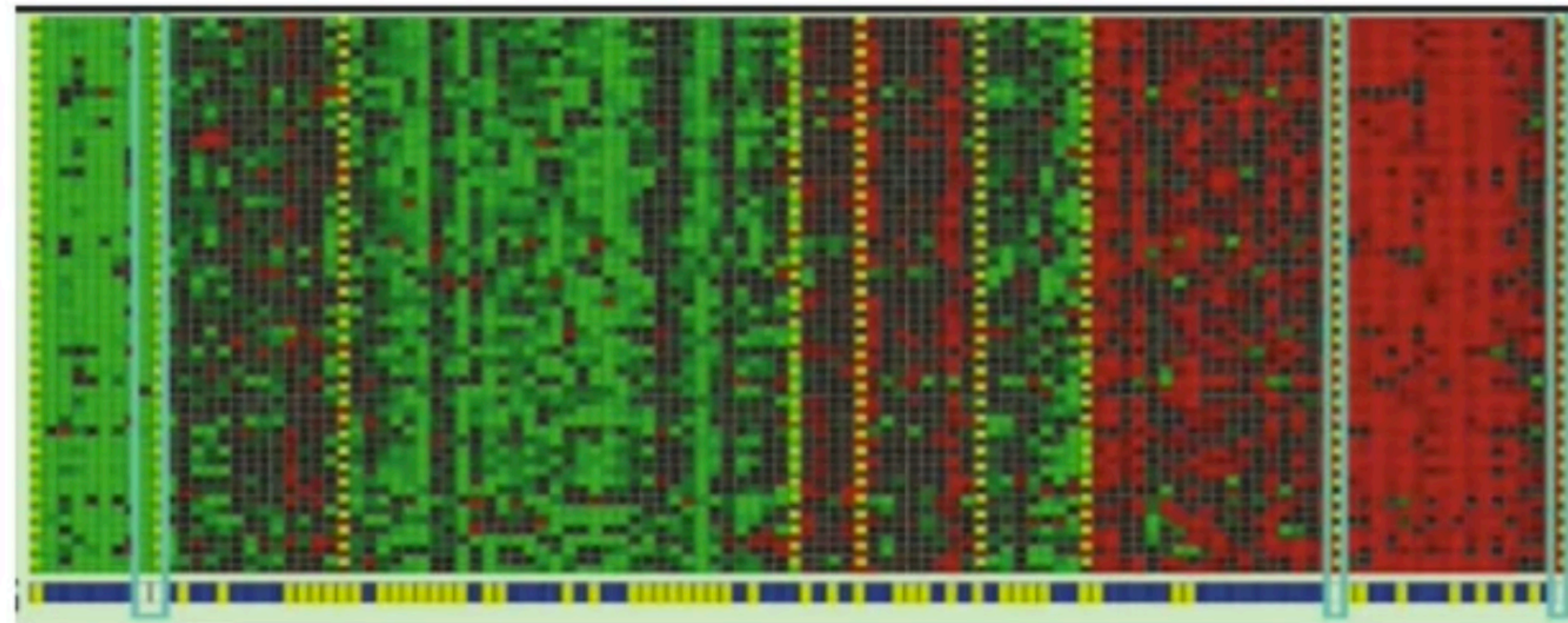
PEOPLE · 6 hours ago

 [View Full Coverage](#)



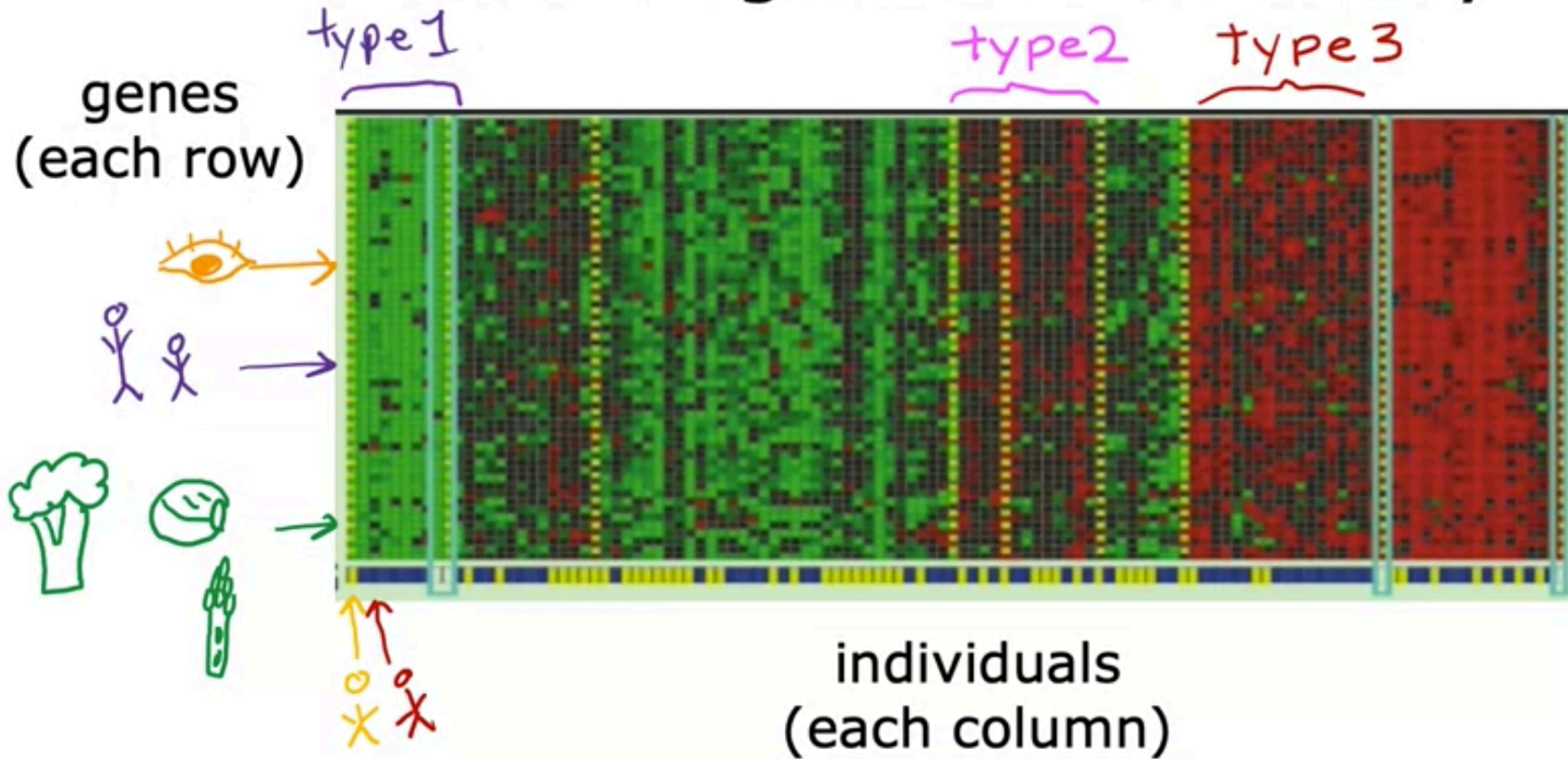
Clustering: DNA microarray

genes
(each row)

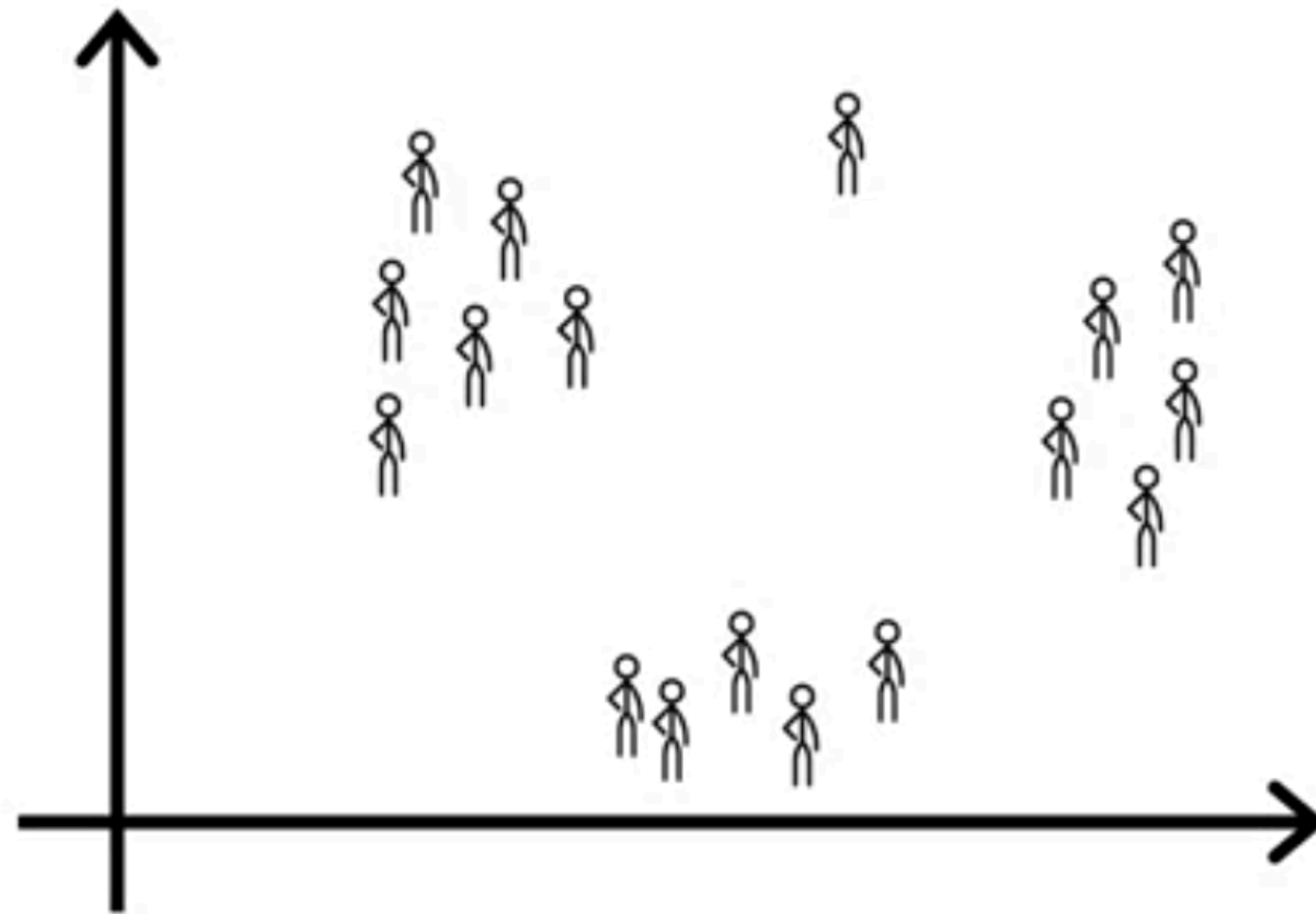


individuals
(each column)

Clustering: DNA microarray

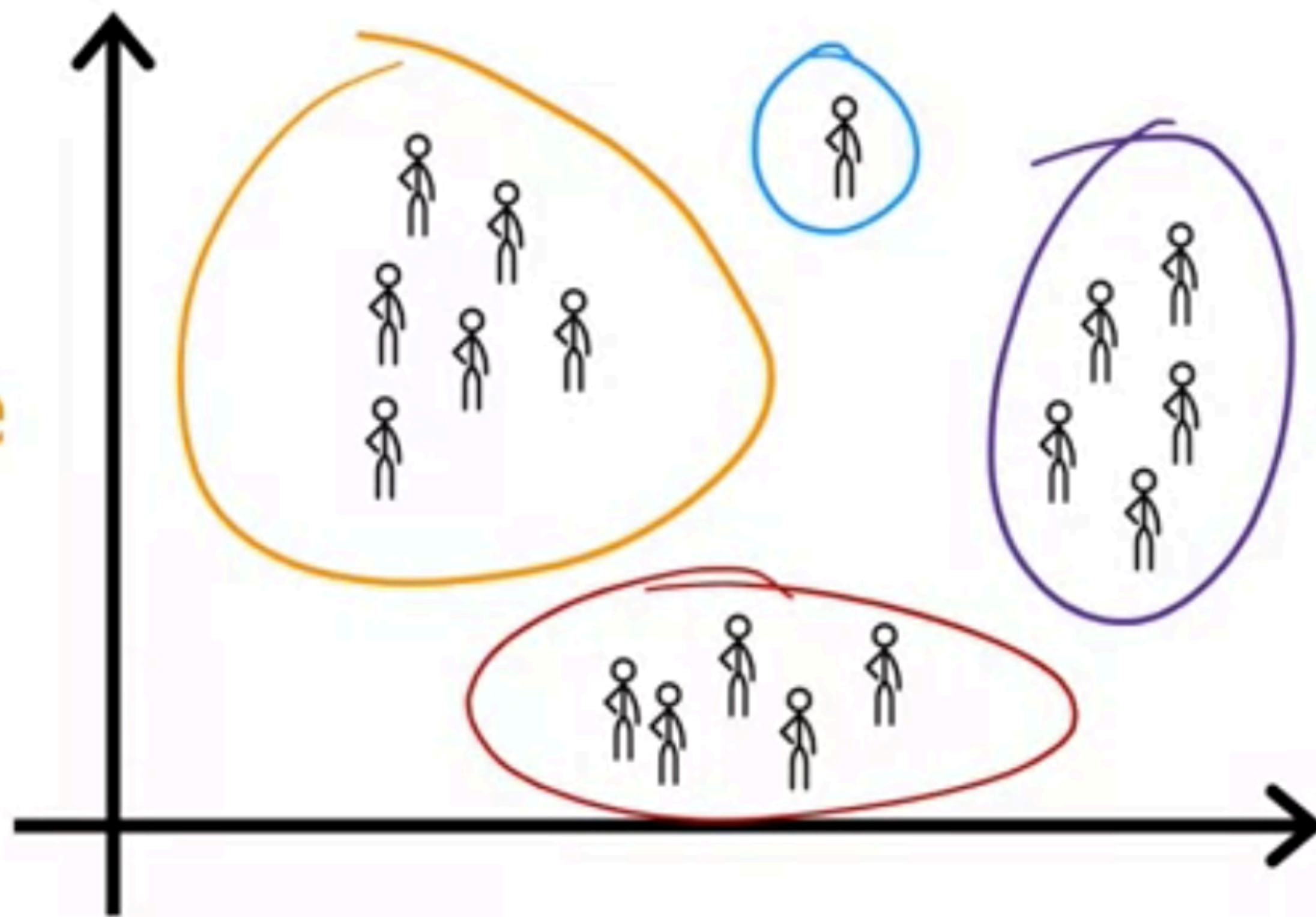


Clustering: Grouping customers



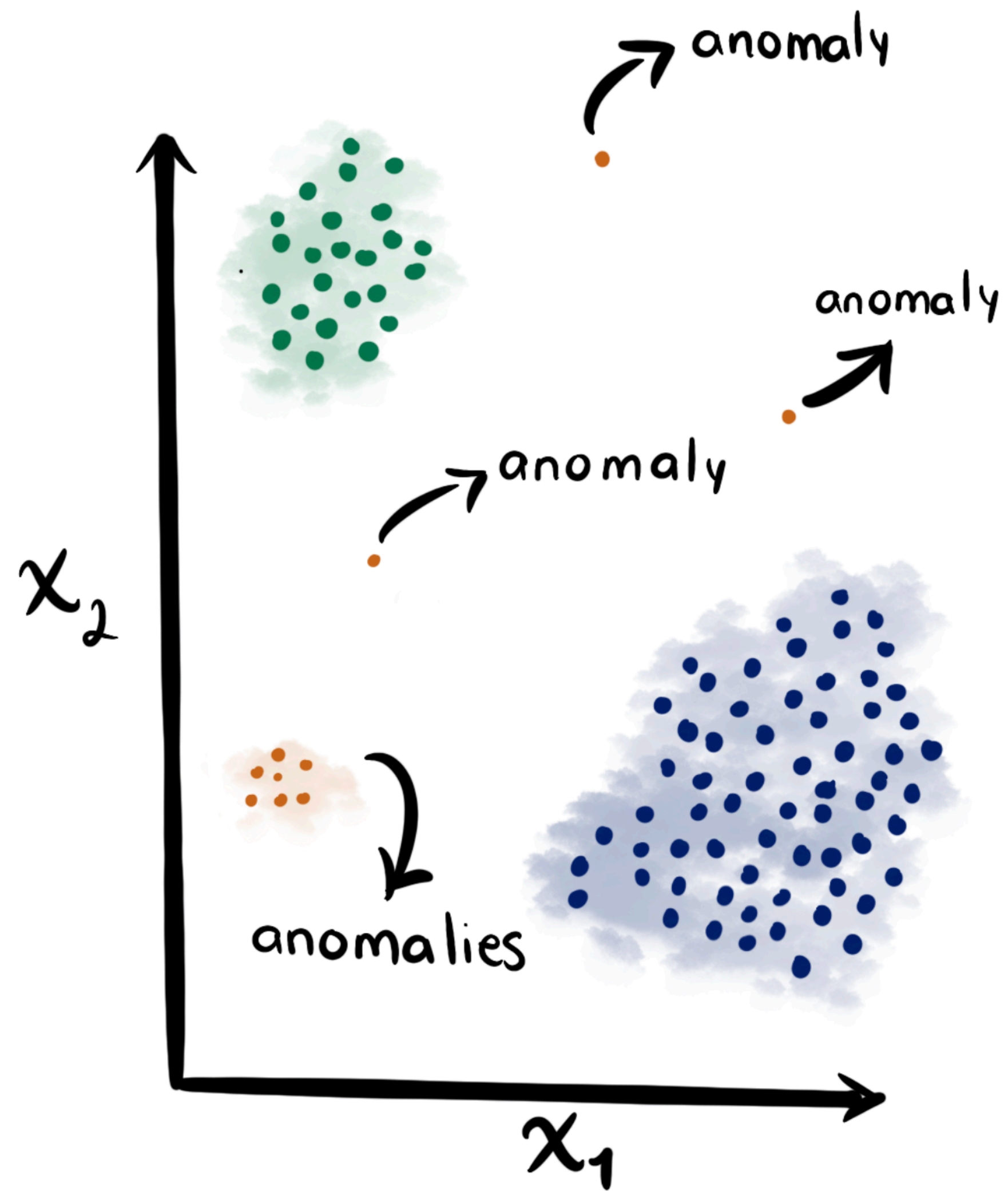
Credit: Andrew Ng, [Machine Learning](#)

GROUPING

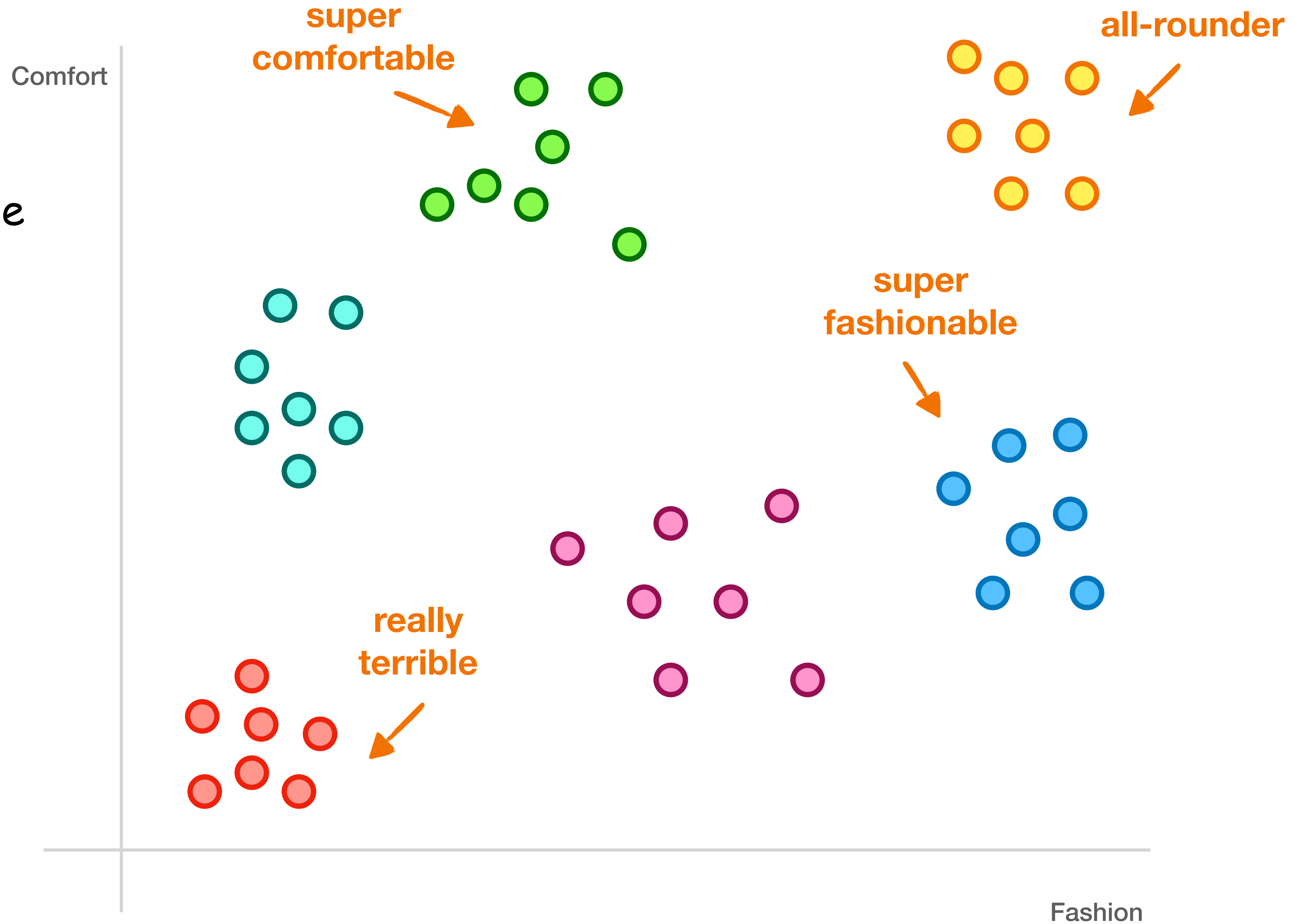


Credit: Andrew Ng, [Machine Learning](#)

ANOMALY D

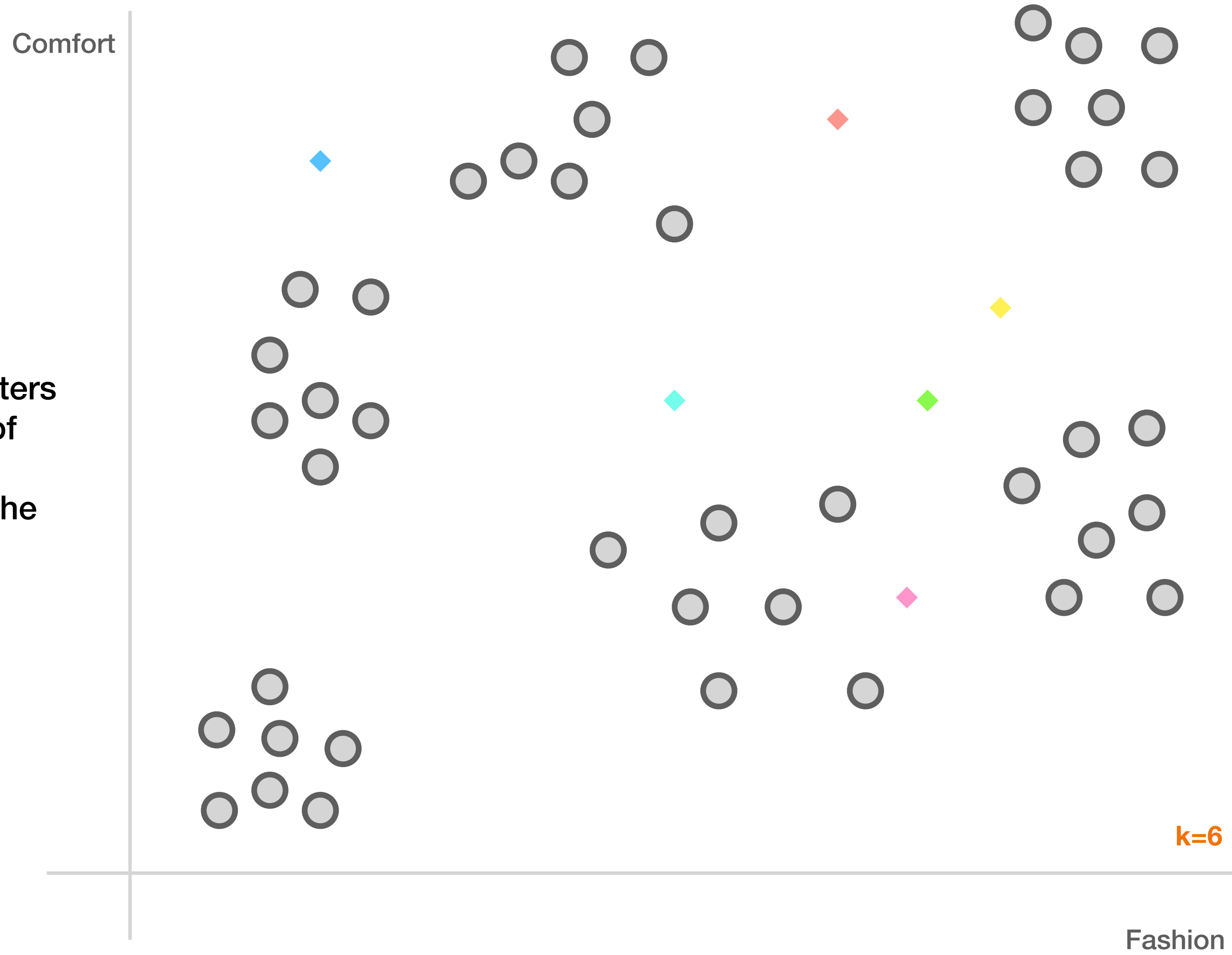


how do we find the clusters?

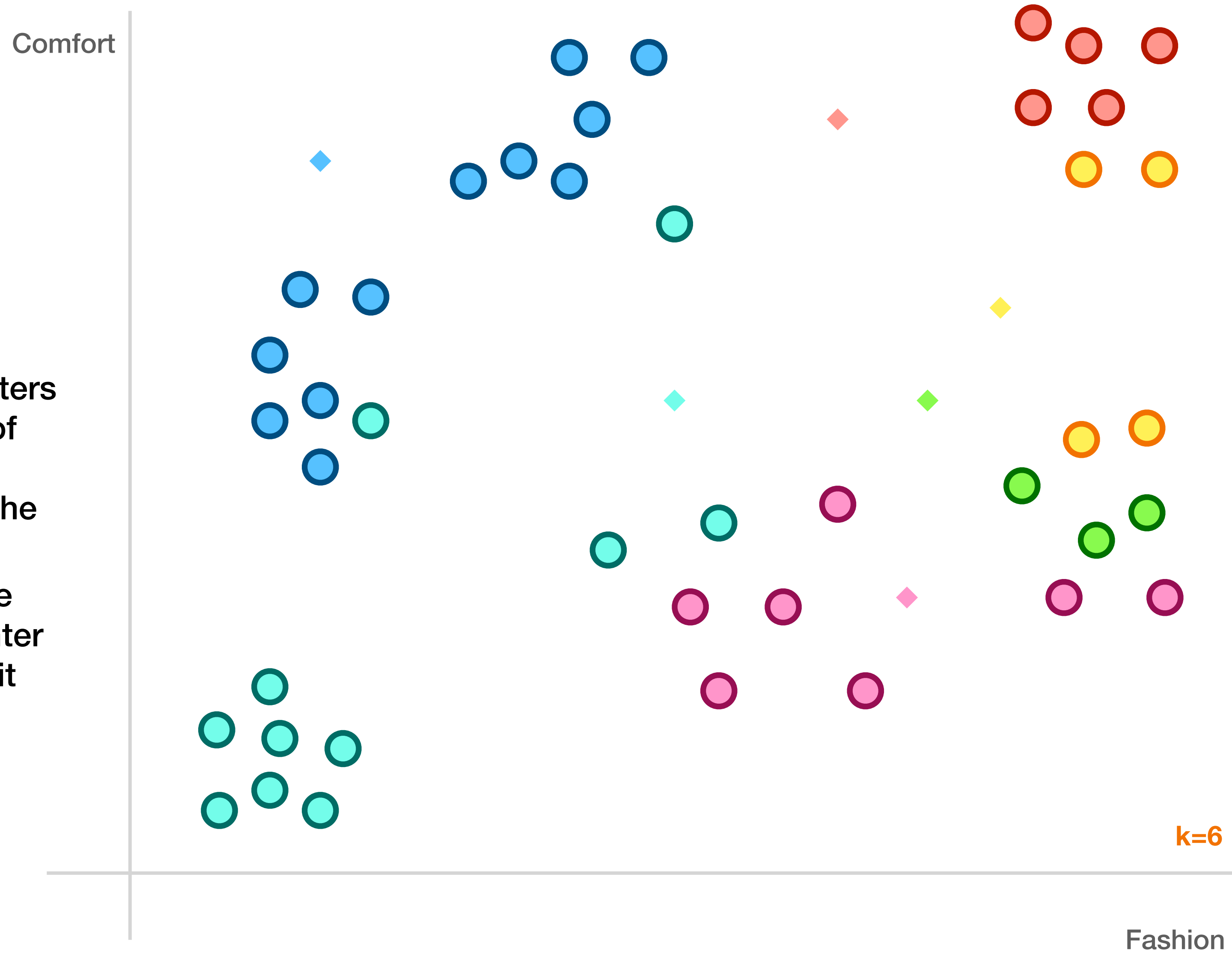


- **clusters** can tell us specifics about the relationship of data
- ...even if they are unlabeled! → **unsupervised learning!**

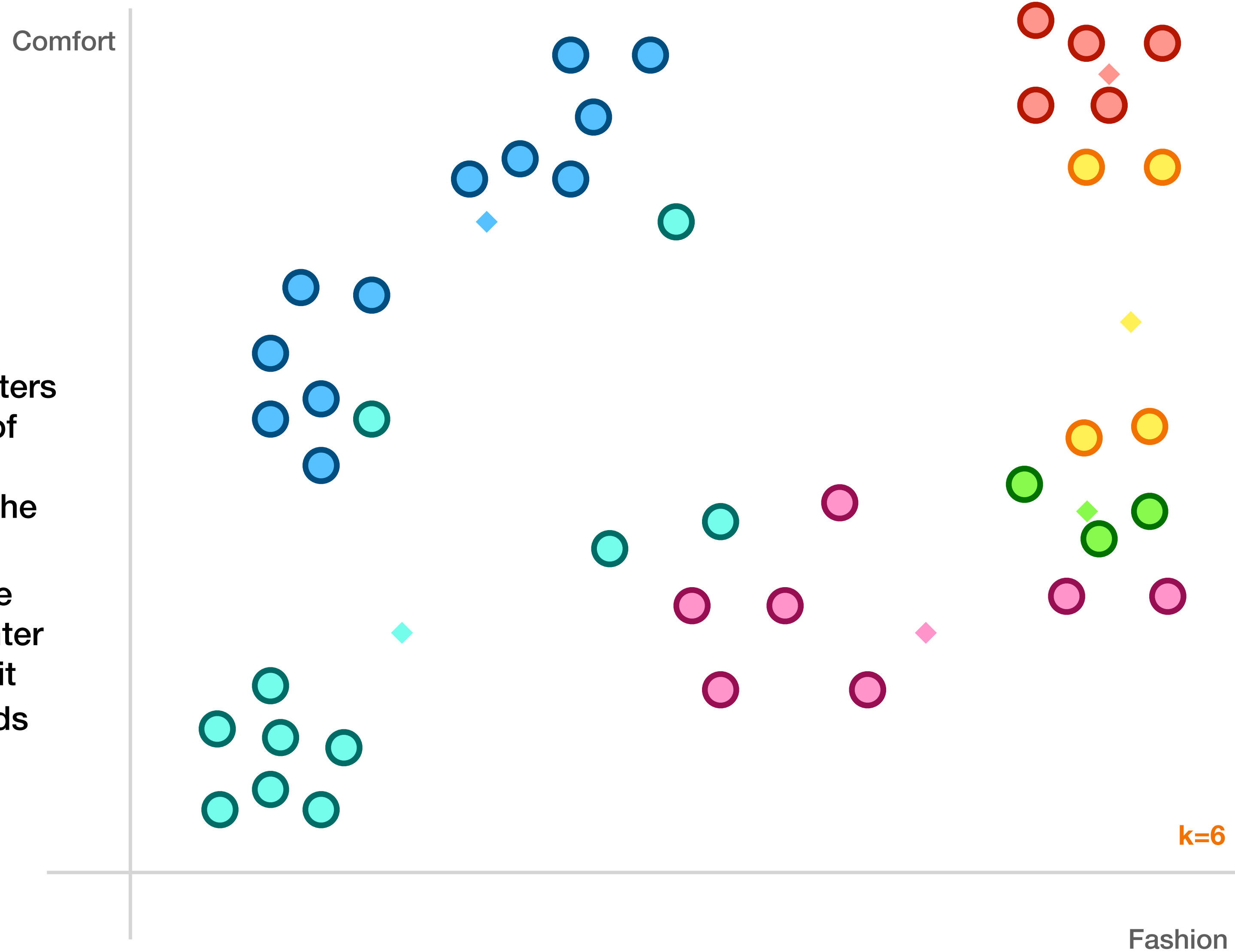
1. pick a K-number of clusters
2. randomly pick a series of "centroids"
3. assign each particle to the centroid closest to it



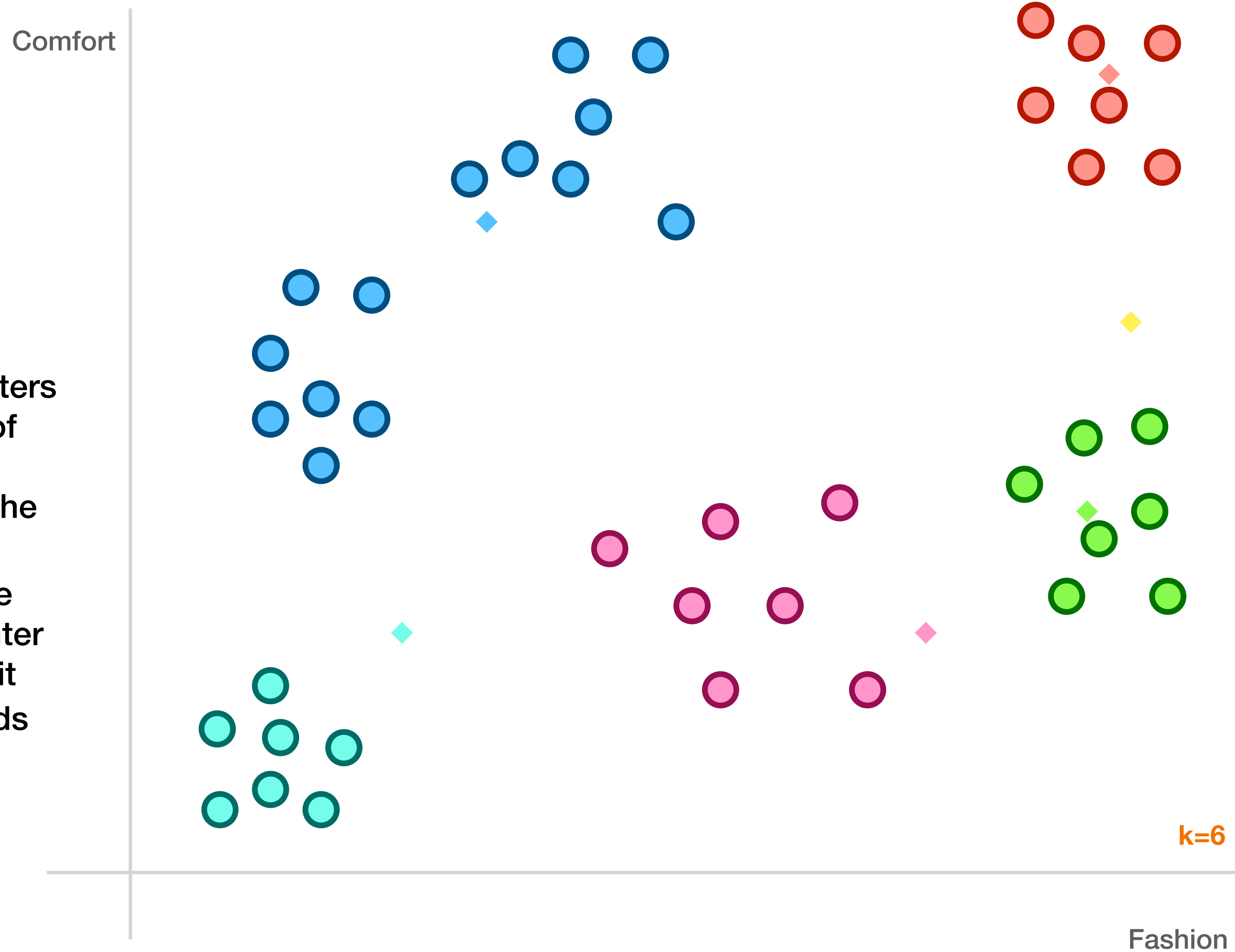
1. pick a K-number of clusters
2. randomly pick a series of “centroids”
3. assign each particle to the centroid closest to it
4. move the centroid to the weighted geometric center of samples assigned to it



1. pick a K-number of clusters
2. randomly pick a series of “centroids”
3. assign each particle to the centroid closest to it
4. move the centroid to the weighted geometric center of samples assigned to it
5. Repeat 3-4 until centroids stop moving!



1. pick a K-number of clusters
2. randomly pick a series of “centroids”
3. assign each particle to the centroid closest to it
4. move the centroid to the weighted geometric center of samples assigned to it
5. Repeat 3-4 until centroids stop moving!

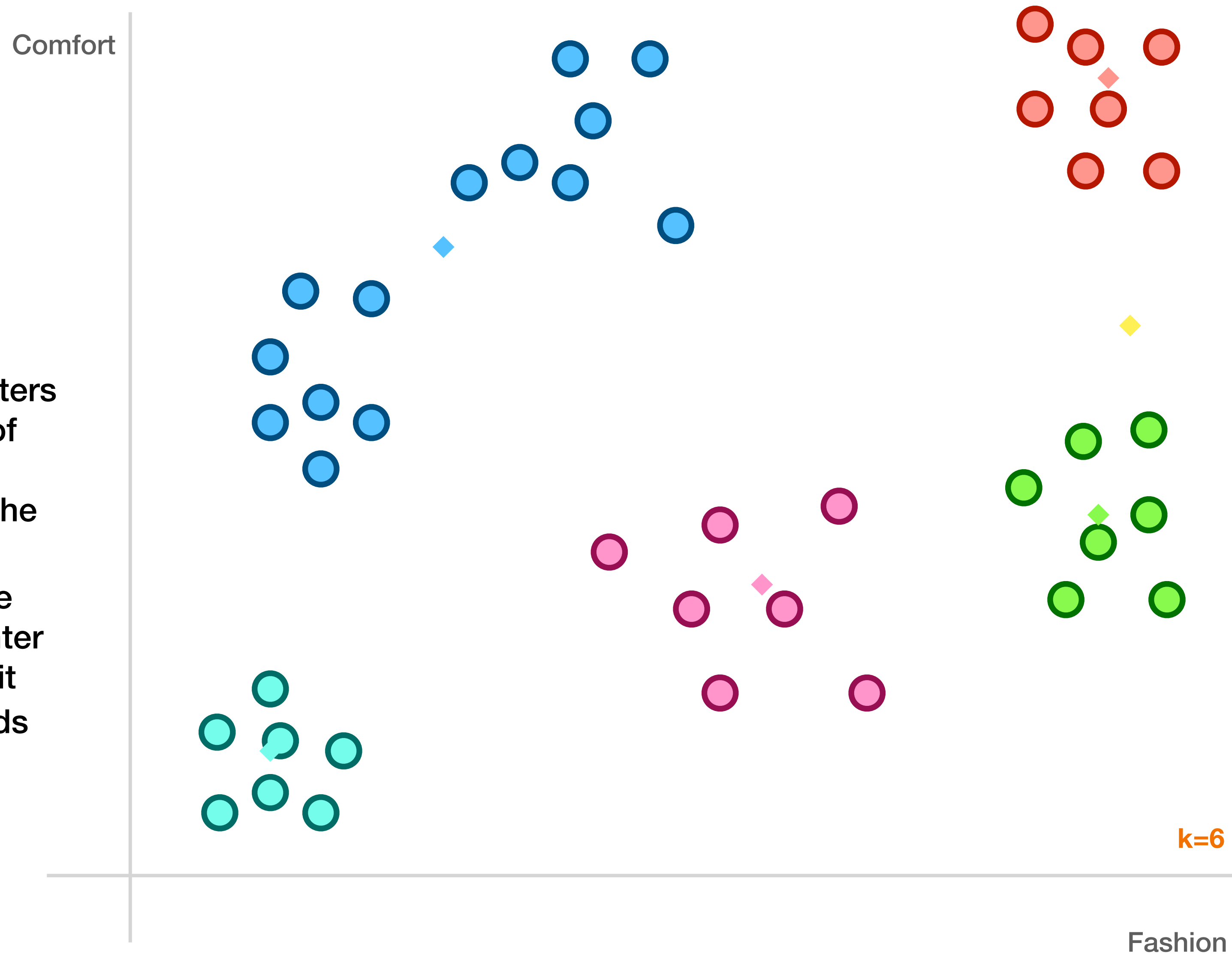


k=6

Fashion

k-means clustering

1. pick a K-number of clusters
2. randomly pick a series of “centroids”
3. assign each particle to the centroid closest to it
4. move the centroid to the weighted geometric center of samples assigned to it
5. Repeat 3-4 until centroids stop moving!



Did we get back
the same clusters?
Nope. And that's OK.



Did we get back the same clusters?

Nope. And that's OK.

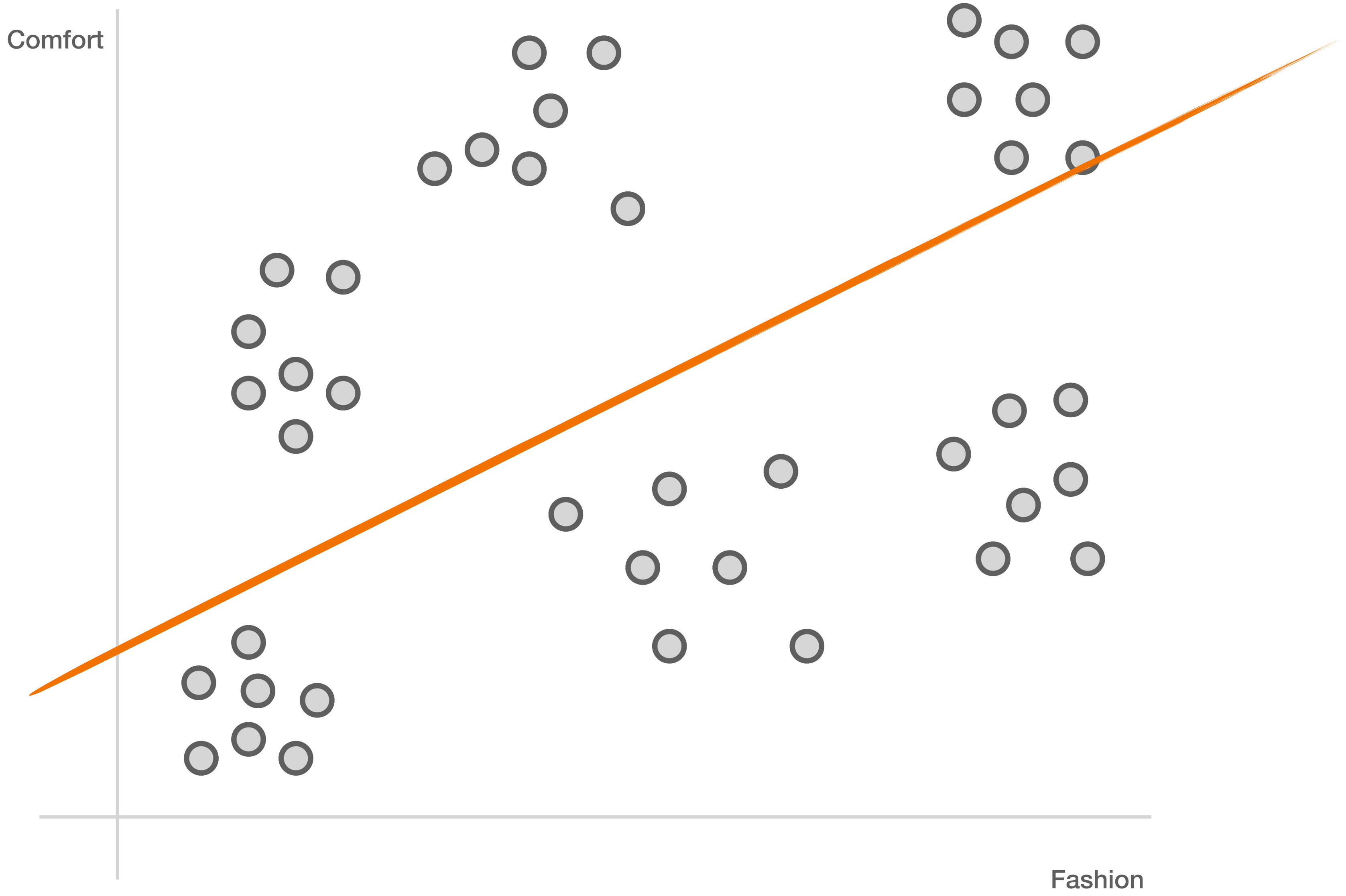
- **K-means** is an *indeterministic* algorithm—it has built-in randomness

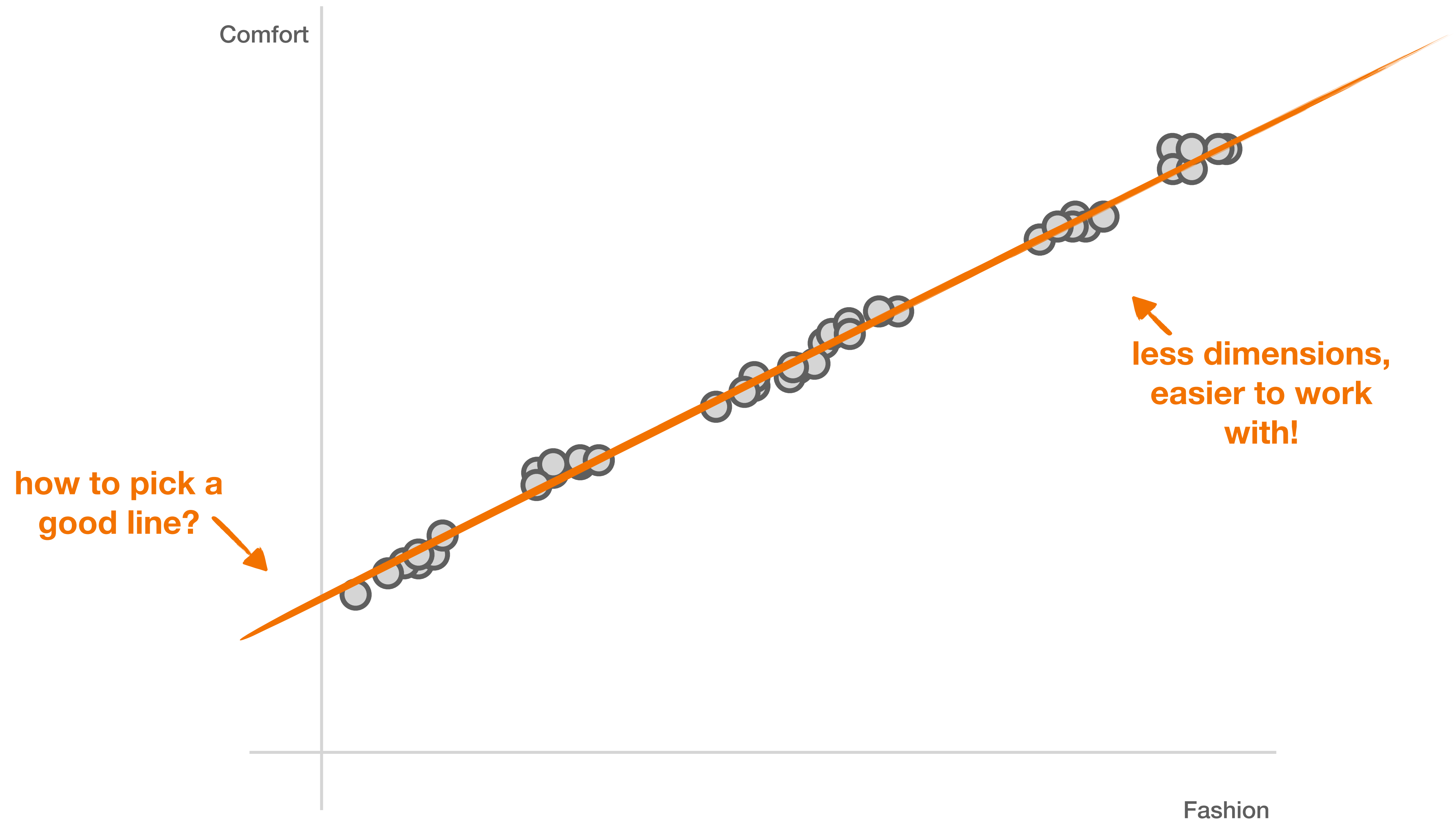
UNSUPERVISED LEARNING

Clustering

Dimension reduction

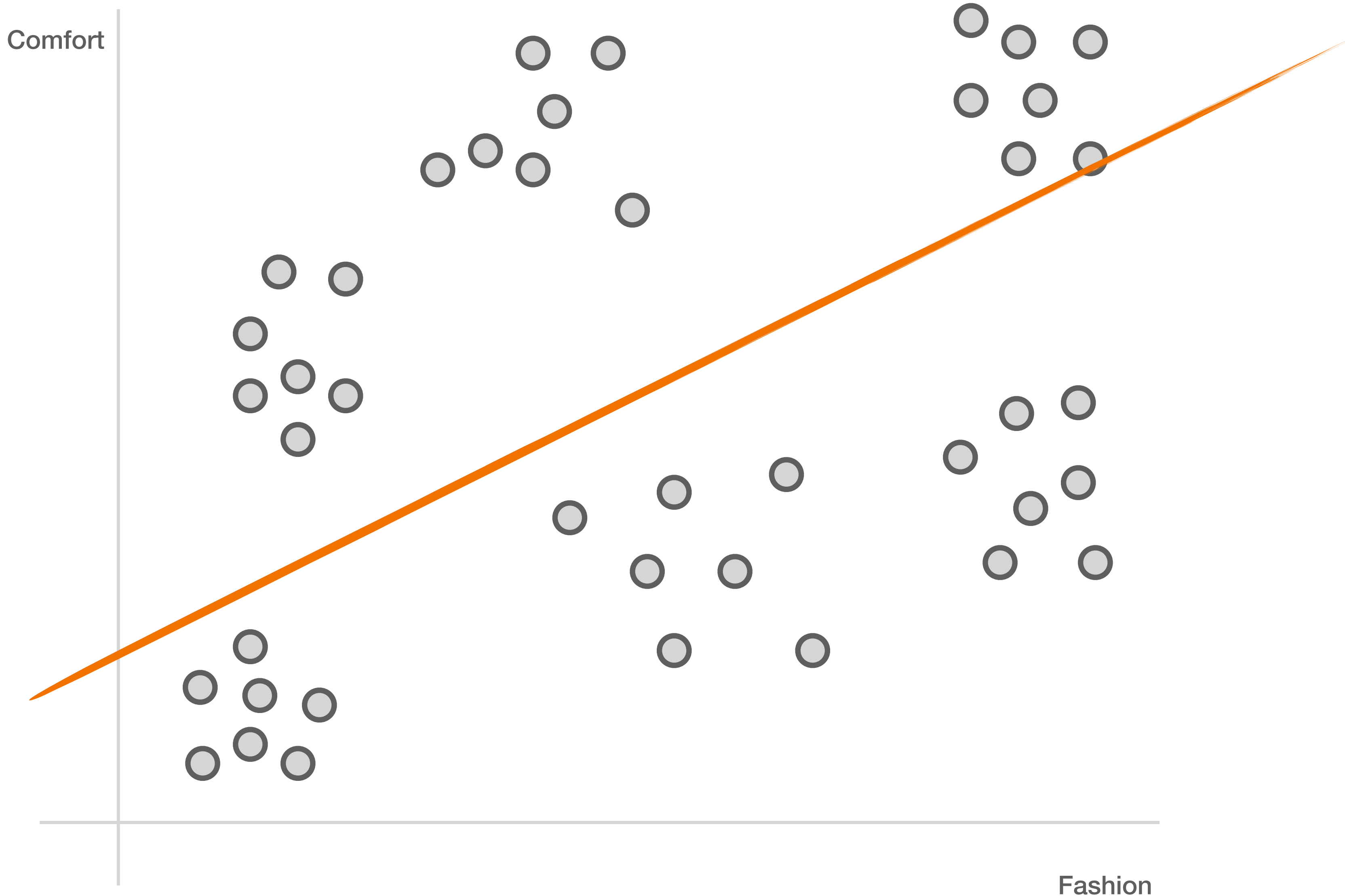
Principle Component Analysis



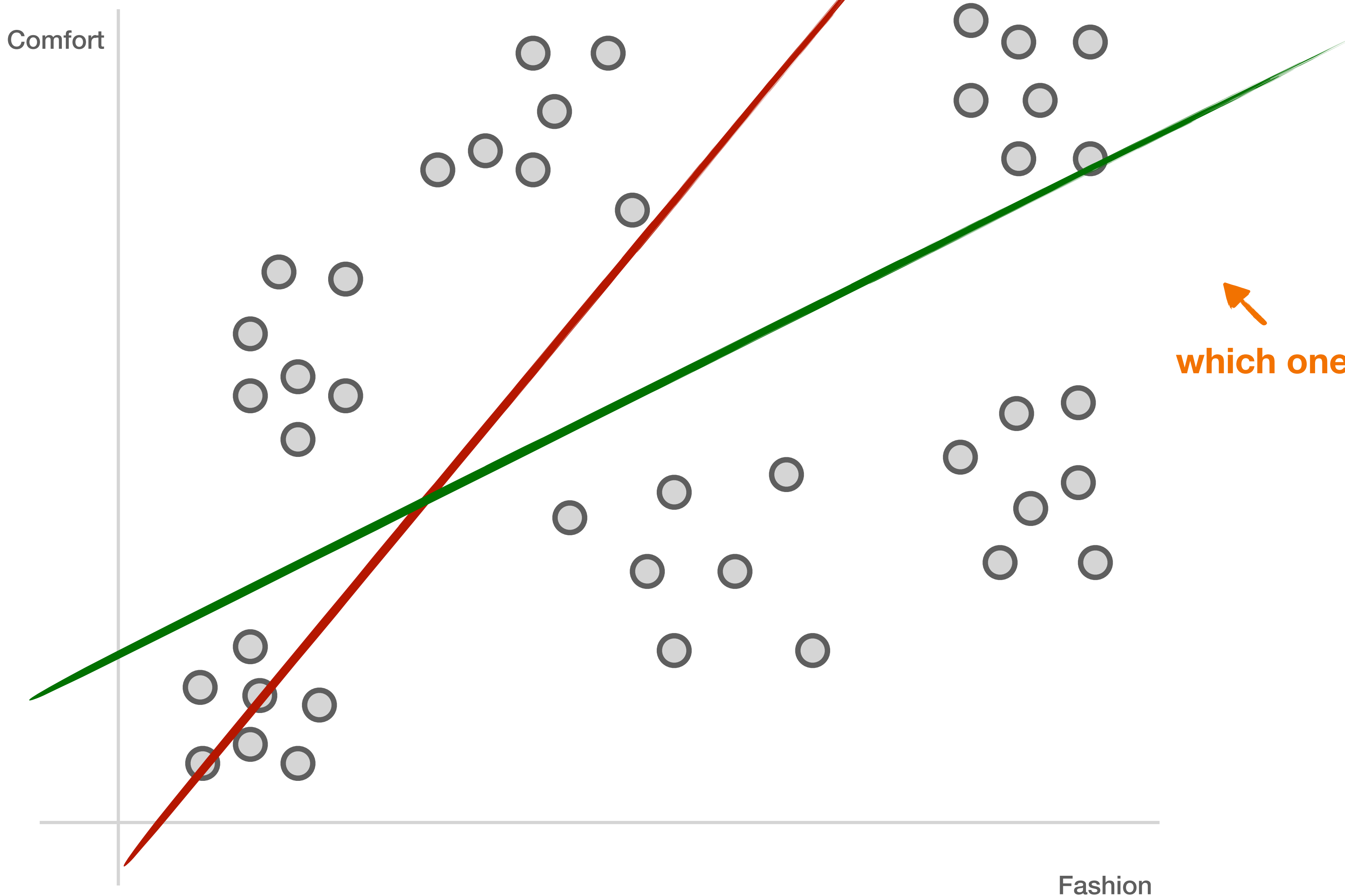


- by **projecting** the samples down to a smaller dimension, they are easier to work with.
(because the centroids have less “space” to move around)

how to pick a good line?

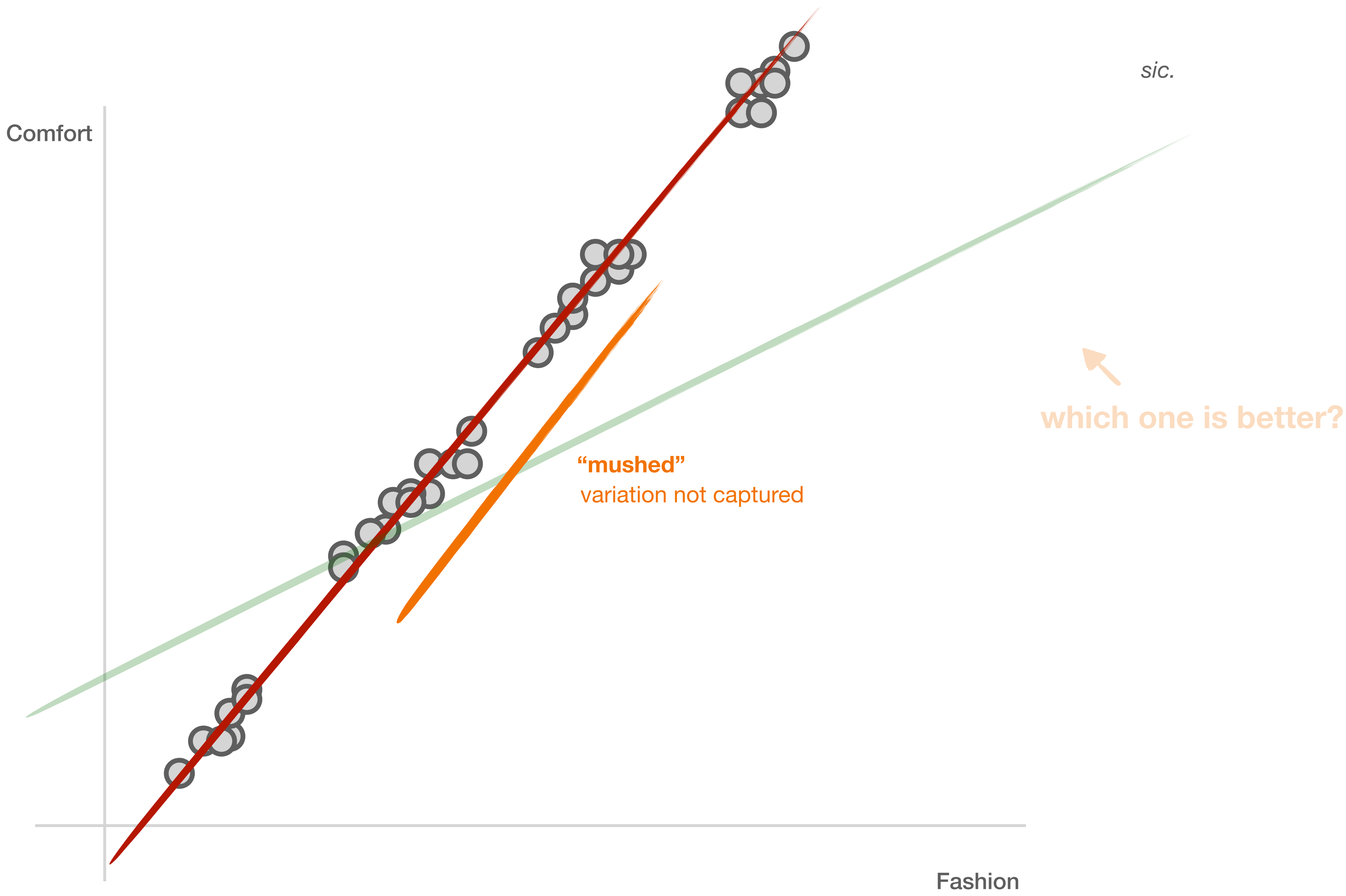


how to pick a good line?



which one is better?

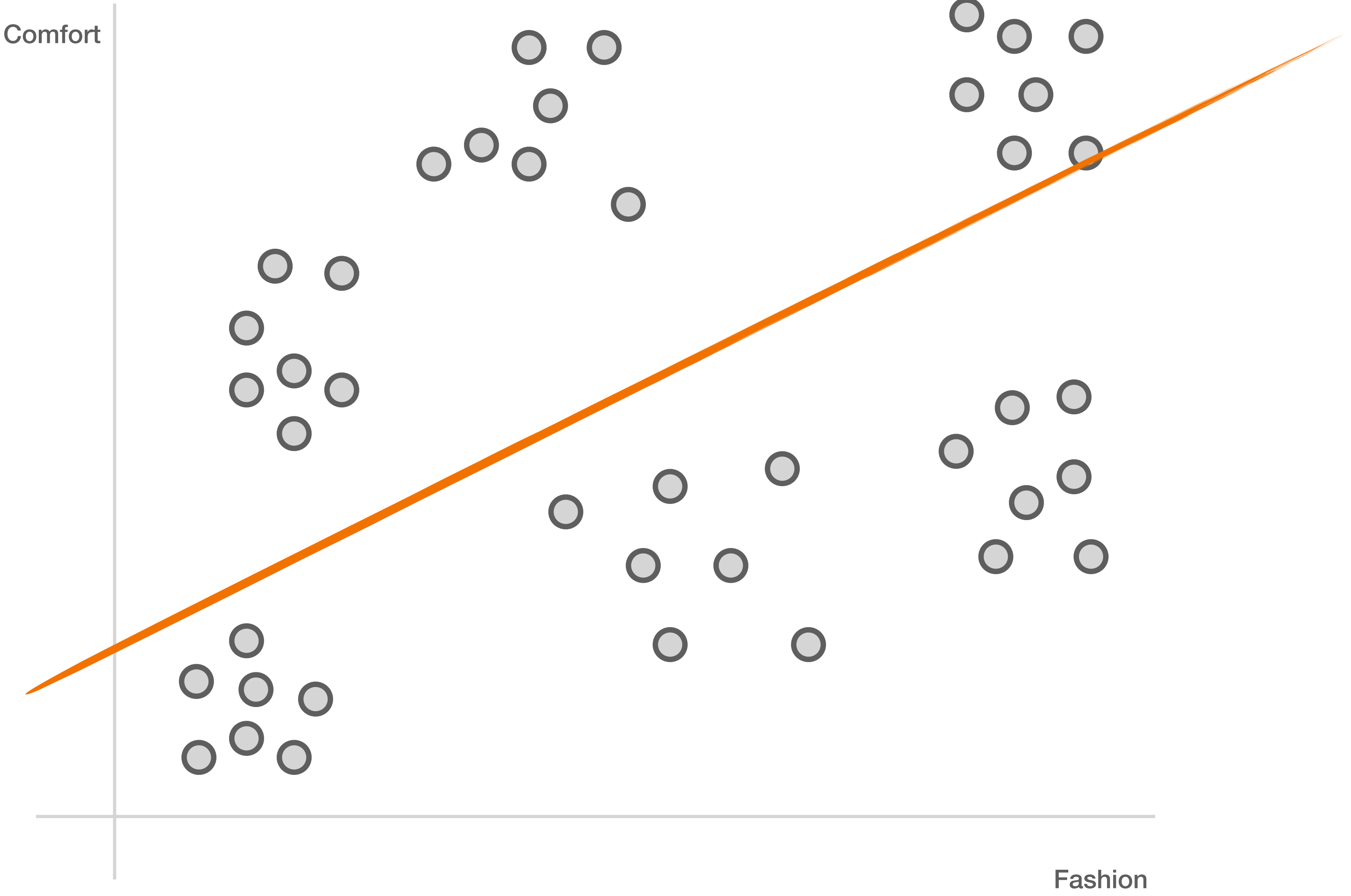
how to pick a good line?



■ a good “**projection**” captures the **variation** in the data

how to pick a good line?

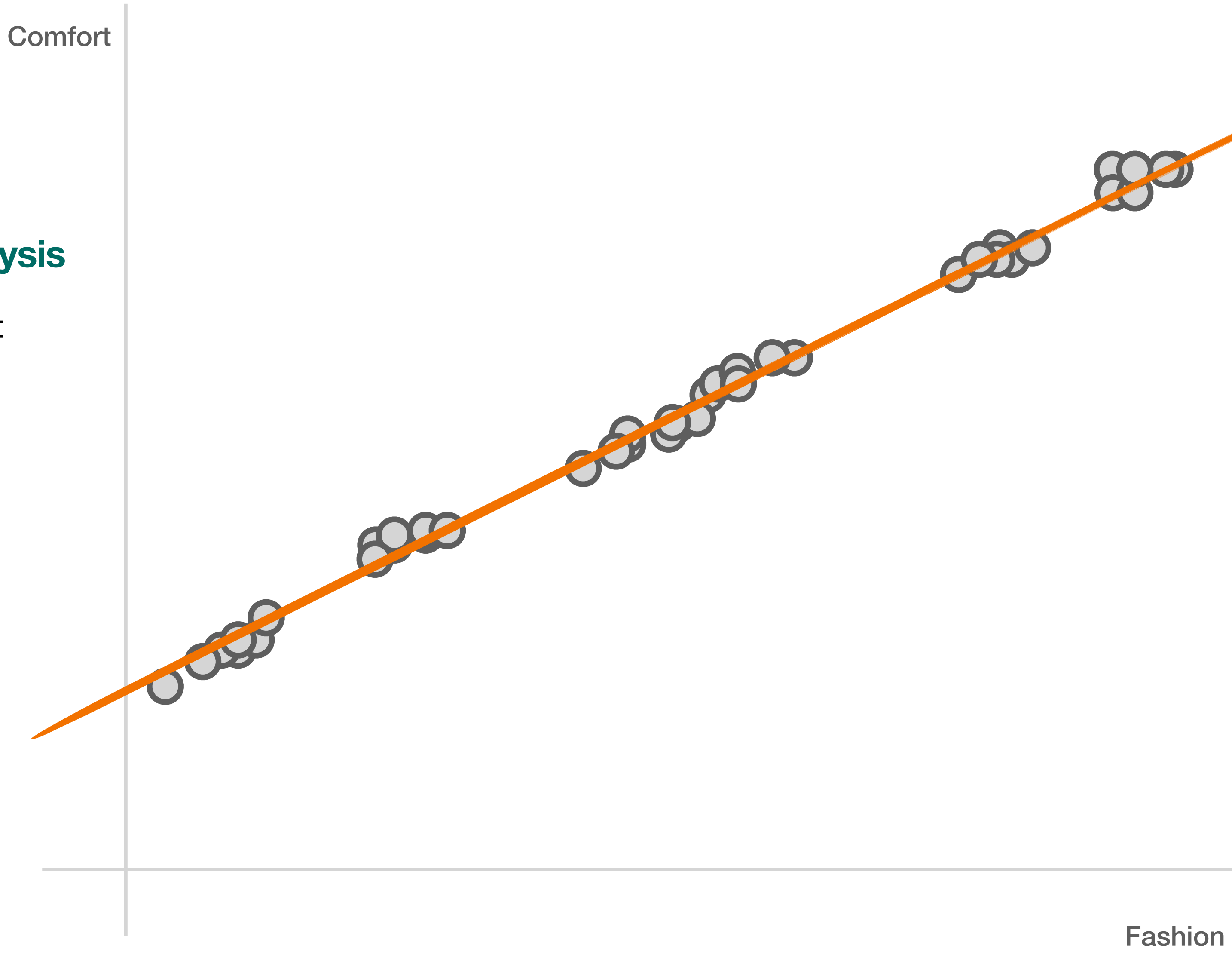
- Find a good line (**basis**) that **maximizes variation**



how to pick a good line?

principle component analysis

- Find a good line (**basis**) that **maximizes variation**
- **Project** samples down



But, can we formalize it?



Principle Component Analysis

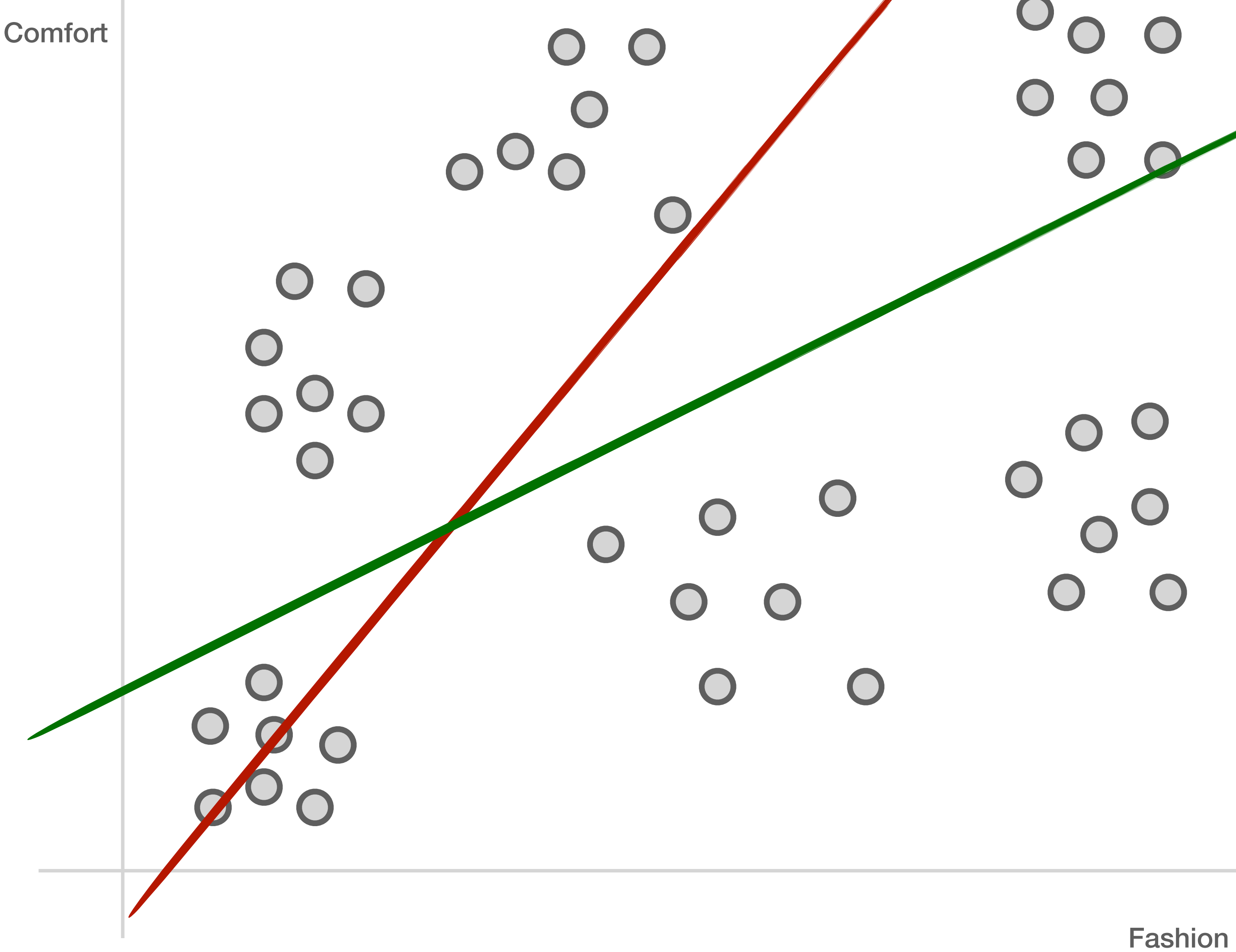
First Component: maximize the explained variance

$$\mathbf{w}_{(1)} = \arg \max_{\|\mathbf{w}\|=1} \left\{ \sum_i (\mathbf{x}_{(i)} \cdot \mathbf{w})^2 \right\}$$

Further Components: maximize the explained variance in remainders

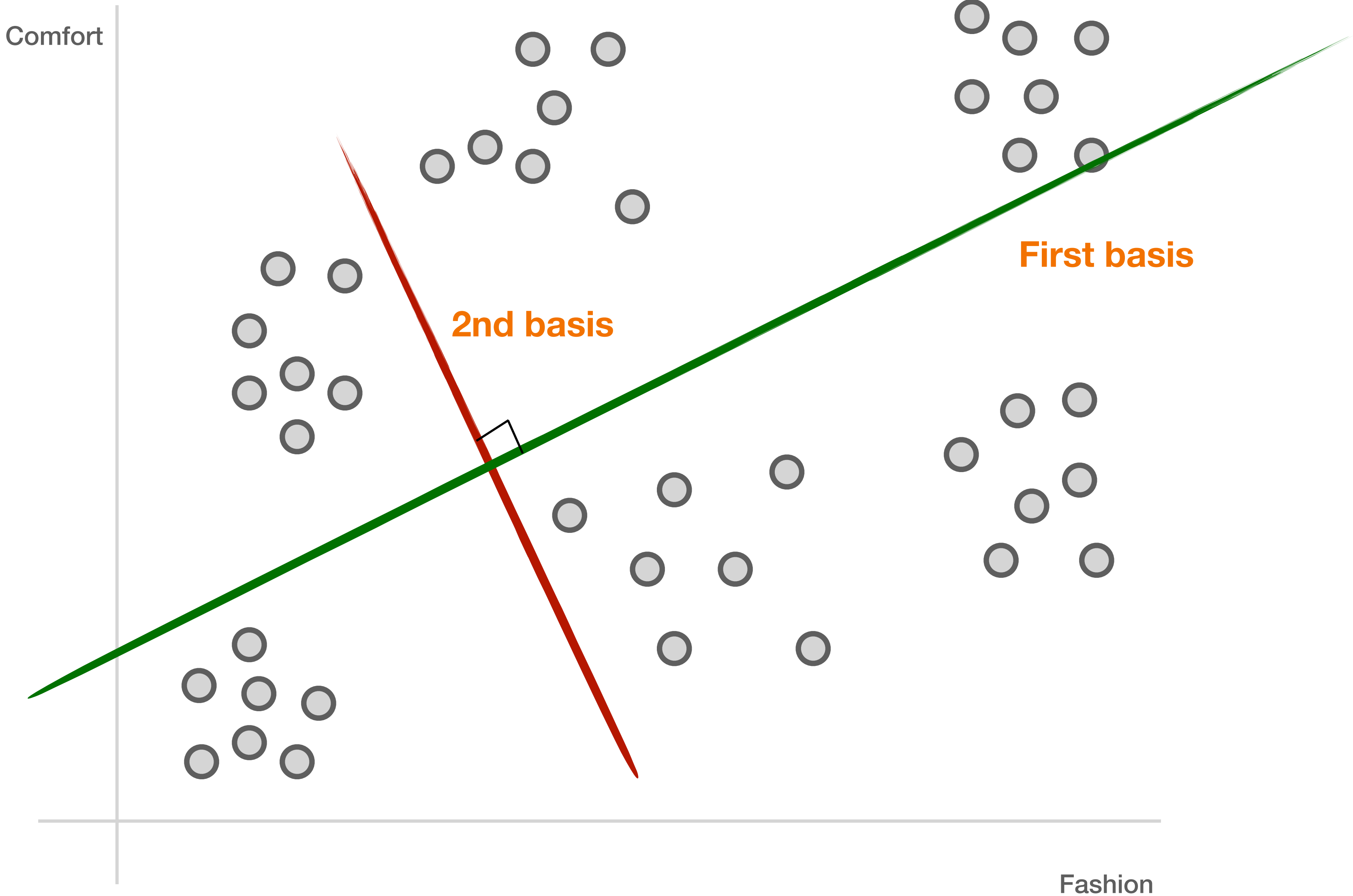
$$\mathbf{w}_{(k)} = \arg \max_{\|\mathbf{w}\|=1} \left\{ \left\| \hat{\mathbf{X}}_k \mathbf{w} \right\|^2 \right\} : \quad \hat{\mathbf{X}}_k = \mathbf{X} - \sum_{s=1}^{k-1} \mathbf{X} \mathbf{w}_{(s)} \mathbf{w}_{(s)}^\top$$

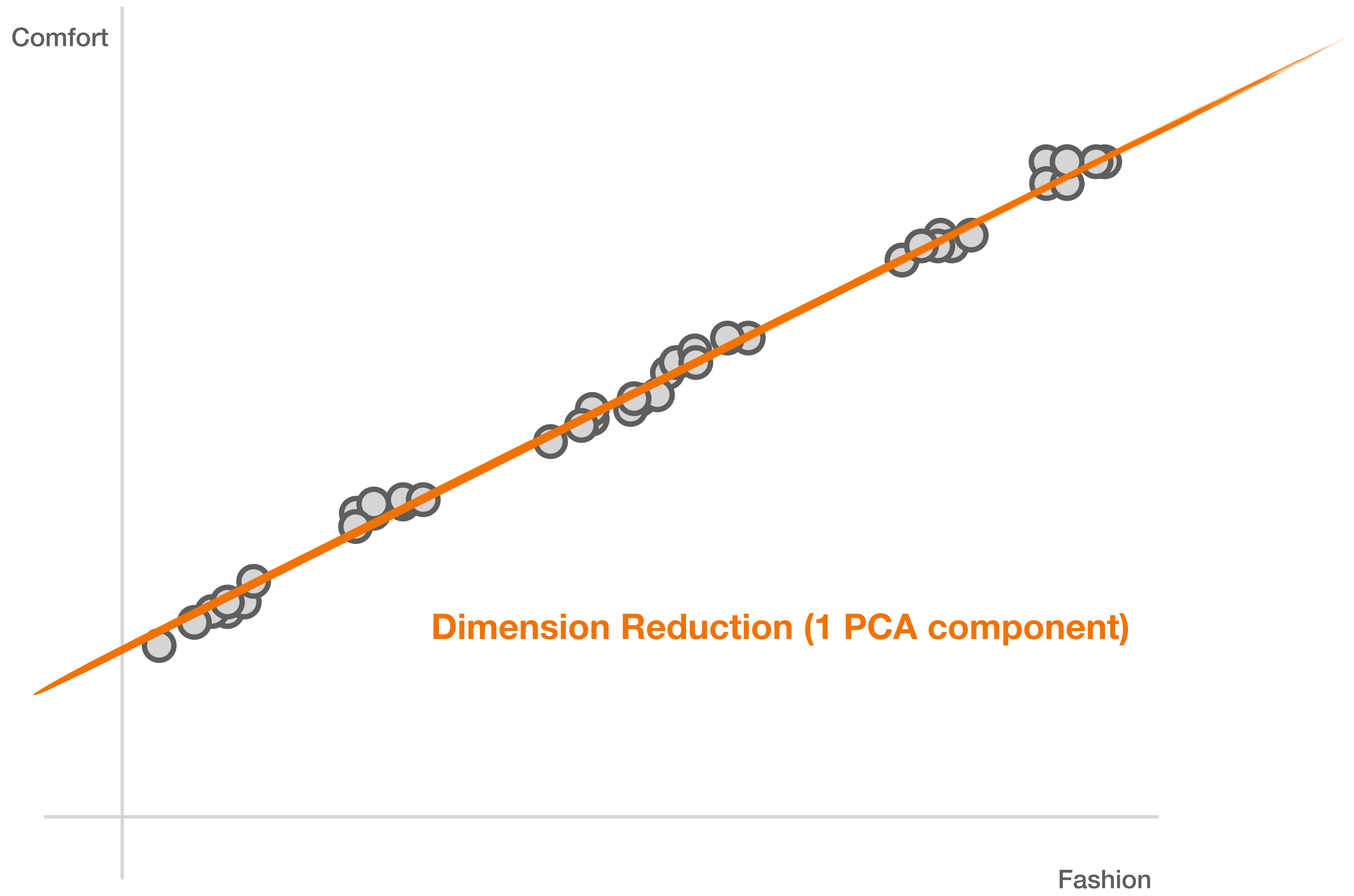
First Component



which one is better?

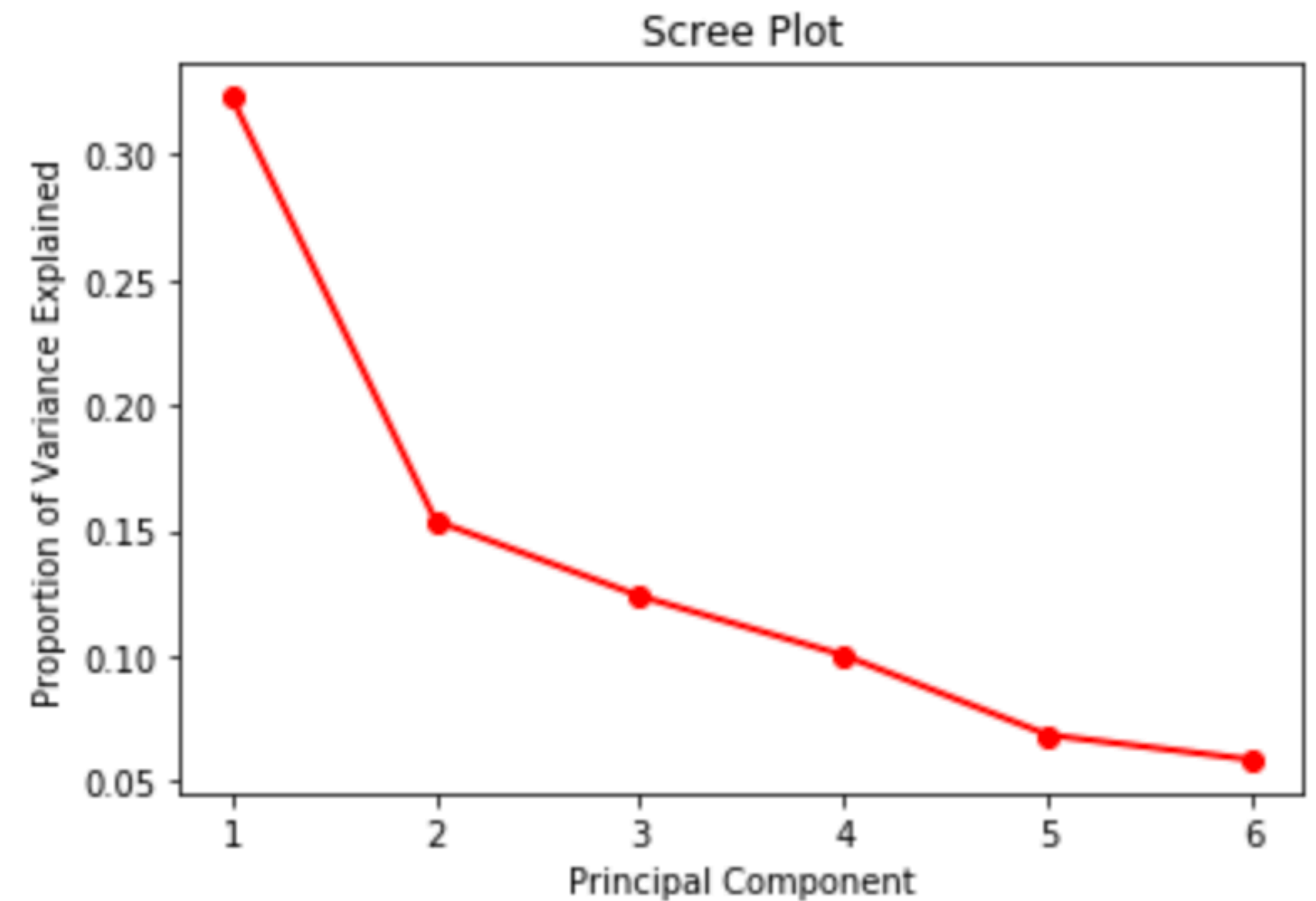
PCA





- Linear transformation of original features
- Dimension reduction
- Compression
- Denoise
- Lose interpretability

- How to choose k (hyperparameter)



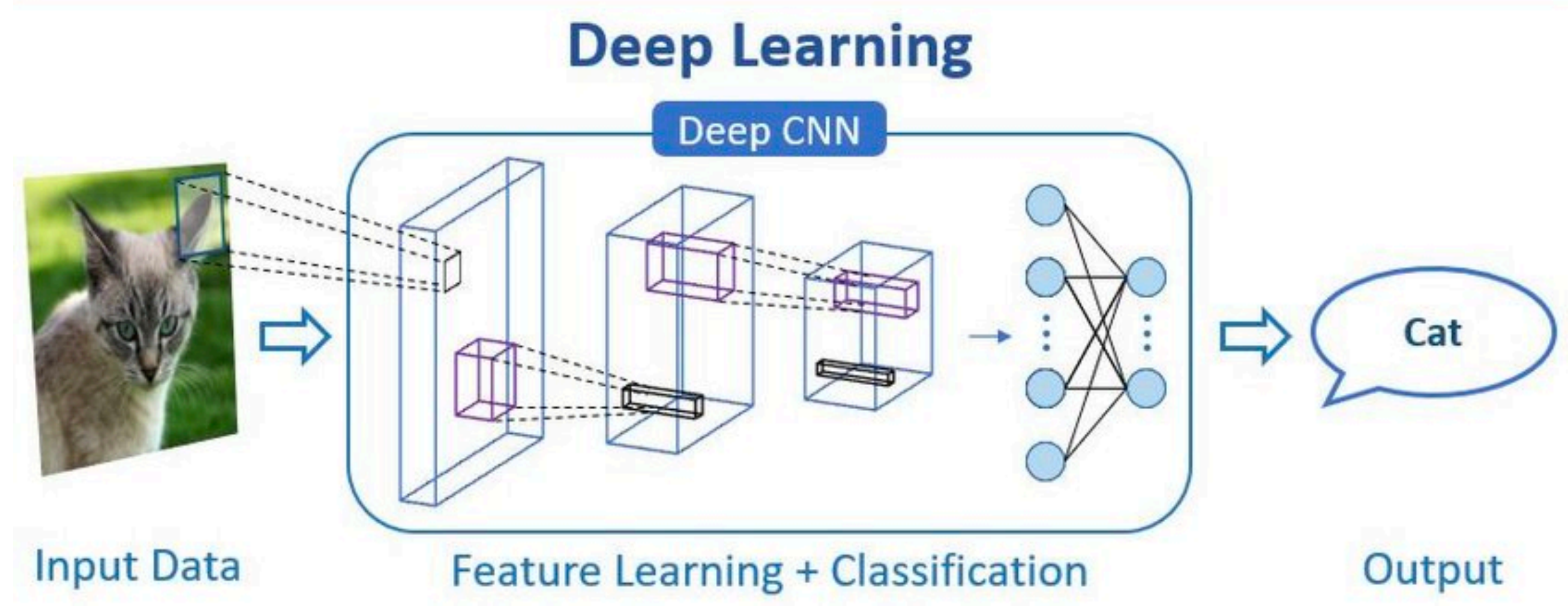
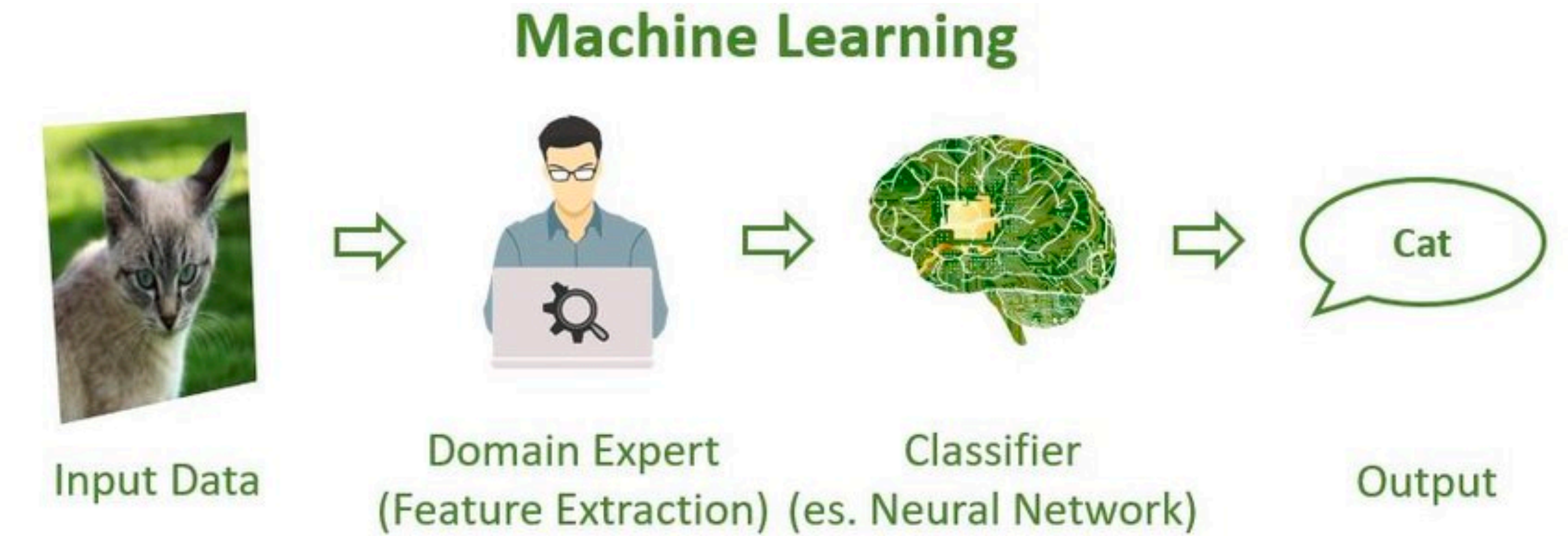


POTENTIAL PITFALLS

Things that can go wrong

- Inconsistent preprocessing (e.g., different scaling/normalization)
- Data leakage (e.g., temporal or mixing subjects)
- Model is used on test data that has changed
- Selecting appropriate metrics (e.g., is 99% accuracy good enough?)
- Hidden confounders (e.g., golf is correlated with heart attacks)
- Spurious correlations (e.g., hospital ID on images)
- Performance on subgroups may be missing

Classic vs. Deep ML



Cool, so...
what's next?



ENSEMBLE

Wisdom of the crowd

- Guessing the weight of a steer (Sir Francis Calton)
- Key components:
 - Base models with **diversity**
 - Infusion algorithms to integrate base models
- Bagging
- Boosting

Bagging

- Build several instances of a classifier using a subset of the original training data
- Aggregate (Averaging) the results
- Reduce variance
- Ex: random forest

- <https://scikit-learn.org/stable/modules/ensemble.html#bagging-meta-estimator>

Boosting

- Fit a sequence of weak learners on repeatedly modified versions of the data.
- A weighted majority vote
- Ex: Adaboost, Gradient boosting

SELF-SUPERVISED LEARNING

Exploiting Unlabeled Data

A lot of unlabeled data is plentiful and cheap, e.g.,

- Document off the web
- Speech samples
- Images and video

But labeling can be expensive

Self-supervised learning



Text Corpus

Nothing is impossible.
Even the word
impossible
says I'm possible



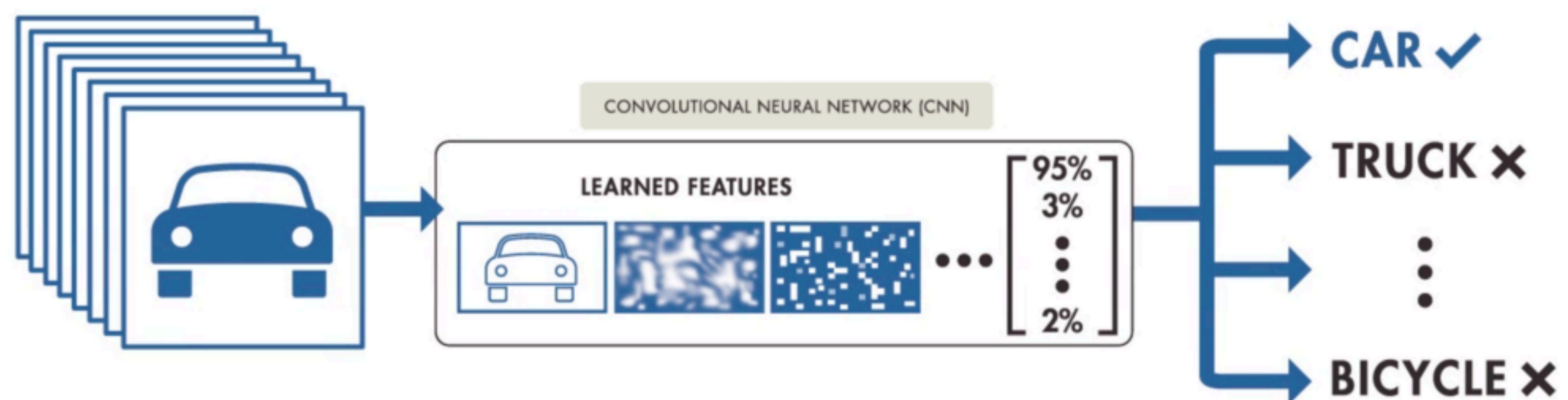
Task: Predict from past

Nothing
Nothing is
Nothing is impossible
...

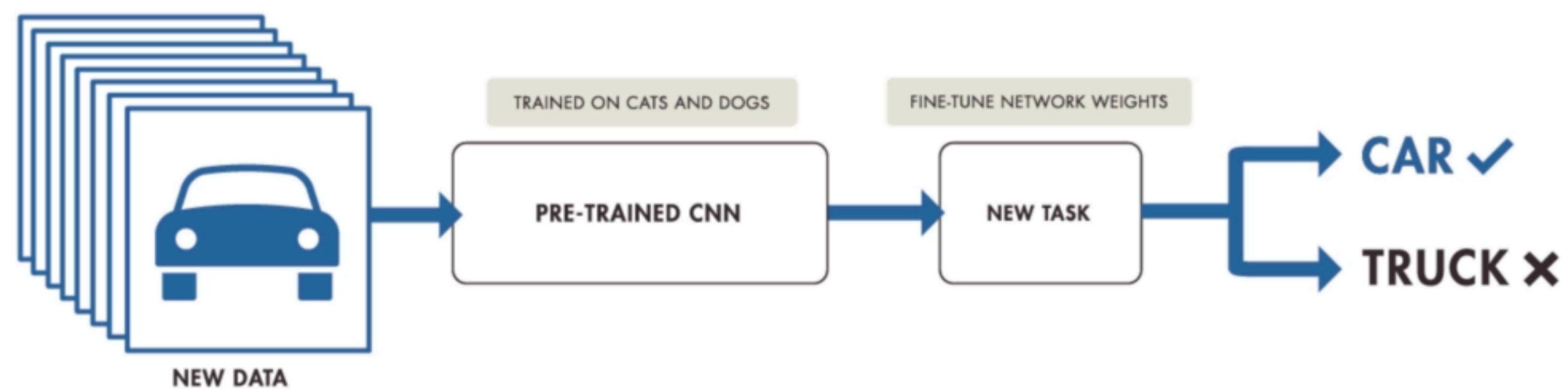
A critical step in training large language models

TRANSFER LEARNING

TRAINING FROM SCRATCH



TRANSFER LEARNING





Pre-training



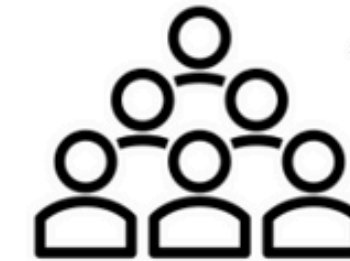
**MLM on
unlabelled data**

word2vec
GloVe
skip-thought
InferSent
ELMo
ULMFIT
GPT
BERT

Fine-tuning



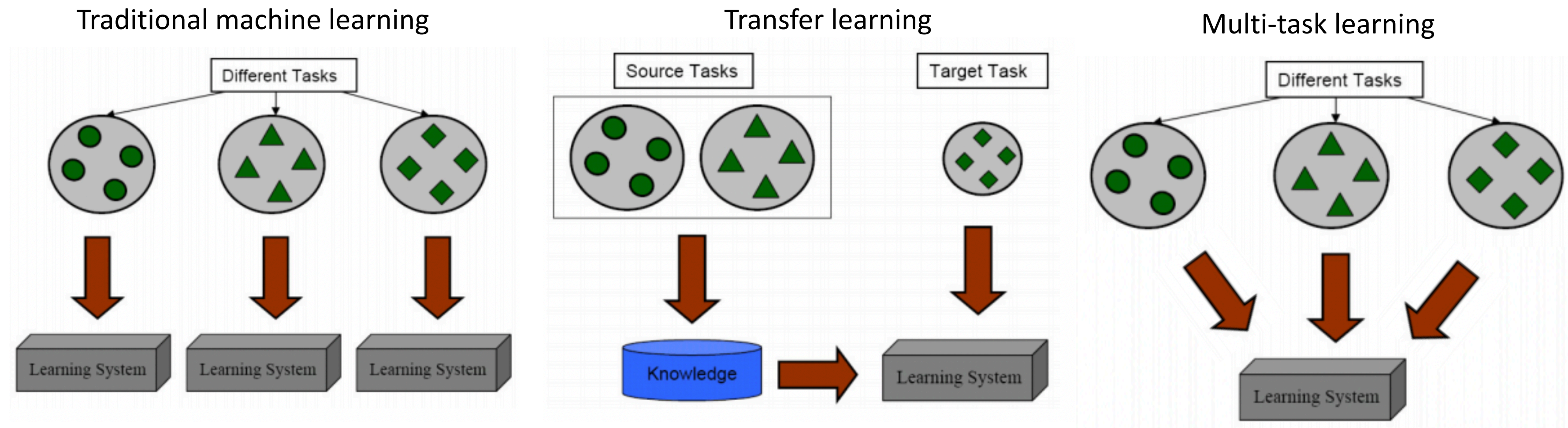
**Cross-entropy
on task labels**



classification
sequence labeling
Q&A

....

Transfer learning and multi-task learning



ACTIVE LEARNING

Active Learning, aka Query Learning

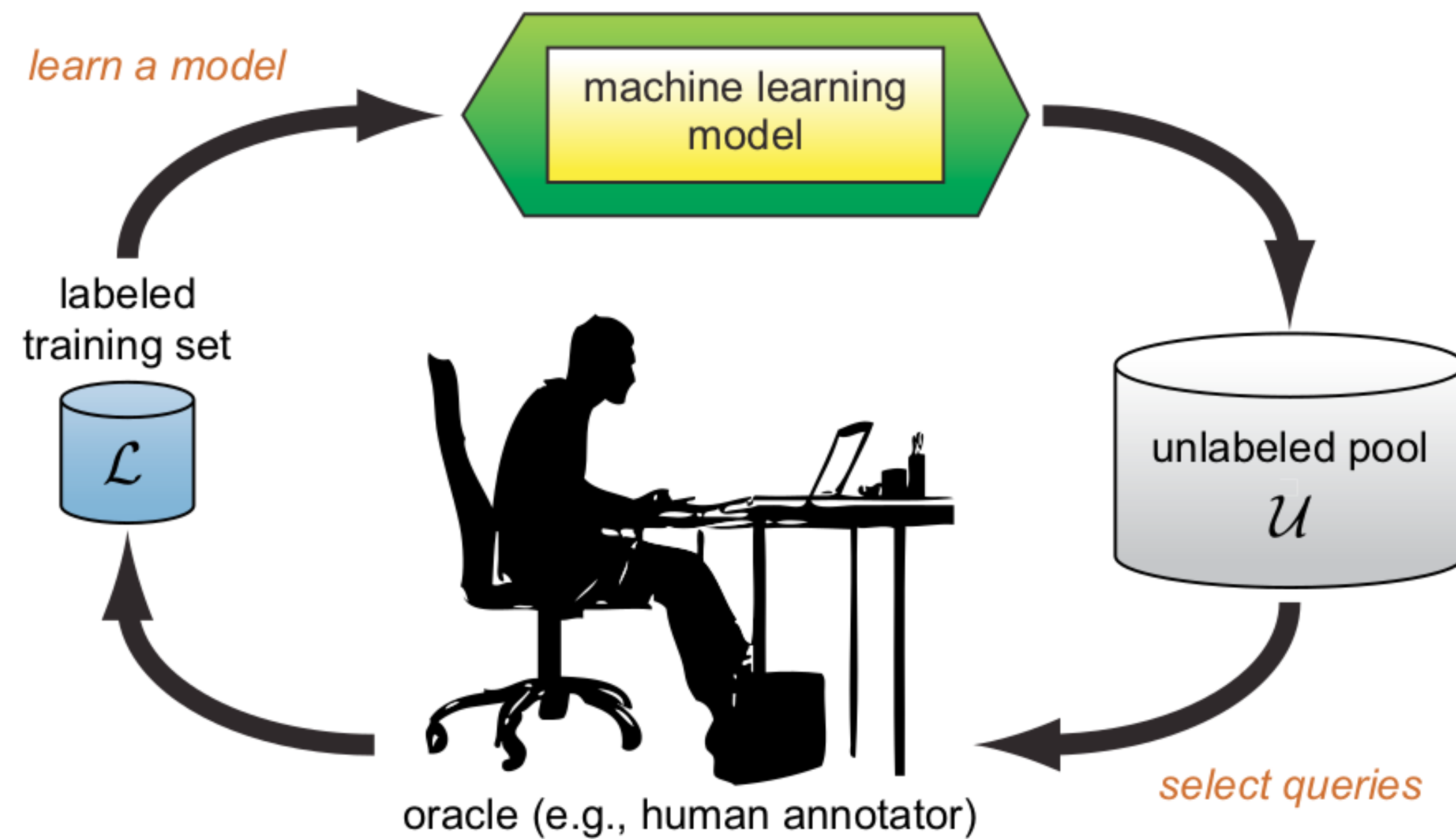


Fig. The active learning cycle.

Repeat

1. Choose unlabeled sample
2. Annotate the chosen unlabeled sample
3. The model trains on the labeled data set

Cheap unlabeled data
Expensive labeled data

Cool, so... what's next?

scikit-learn Install User Guide API Examples Community More ▾

Prev Up Next

1.3. Kernel ridge regression

Other versions

Please cite us if you use the software.

1.4. Support Vector Machines

1.4.1. Classification

1.4.2. Regression

1.4.3. Density estimation, novelty detection

1.4.4. Complexity

1.4.5. Tips on Practical Use

1.4.6. Kernel functions

1.4.7. Mathematical formulation

1.4.8. Implementation details

1.4. Support Vector Machines

Support vector machines (SVMs) are a set of supervised learning methods used for [classification](#), [regression](#) and [outliers detection](#).

The advantages of support vector machines are:

- Effective in high dimensional spaces.
- Still effective in cases where number of dimensions is greater than the number of samples.
- Uses a subset of training points in the decision function (called support vectors), so it is also memory efficient.
- Versatile: different [Kernel functions](#) can be specified for the decision function. Common kernels are provided, but it is also possible to specify custom kernels.

The disadvantages of support vector machines include:

- If the number of features is much greater than the number of samples, avoid over-fitting in choosing [Kernel functions](#) and regularization term is crucial.
- SVMs do not directly provide probability estimates, these are calculated using an expensive five-fold cross-validation (see [Scores and probabilities](#), below).

The support vector machines in scikit-learn support both dense (`numpy.ndarray` and convertible to that by `numpy.asarray`) and sparse (any `scipy.sparse`) sample vectors as input. However, to use an SVM to make predictions for sparse data, it must have been fit on such data. For optimal performance, use C-ordered `numpy.ndarray` (dense) or `scipy.sparse.csr_matrix` (sparse) with `dtype=float64`.

1.4.1. Classification

`SVC`, `NuSVC` and `LinearSVC` are classes capable of performing binary and multi-class classification on a dataset.

Toggle Menu

scikit-learn

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

1.4. Support Vector Machines

1.4.1. Classification

1.4.2. Regression

1.4.3. Density estimation, novelty detection

1.4.4. Complexity

1.4.5. Tips on Practical Use

1.4.6. Kernel functions

1.4.7. Mathematical formulation

1.4.8. Implementation details

Toggle Menu

`SVC` and `NuSVC` are similar methods, but accept slightly different sets of parameters and have different mathematical formulations (see section [Mathematical formulation](#)). On the other hand, `LinearSVC` is another (faster) implementation of Support Vector Classification for the case of a linear kernel. Note that `LinearSVC` does not accept parameter `kernel`, as this is assumed to be linear. It also lacks some of the attributes of `SVC` and `NuSVC`, like `support_`.

As other classifiers, `SVC`, `NuSVC` and `LinearSVC` take as input two arrays: an array `X` of shape `(n_samples, n_features)` holding the training samples, and an array `y` of class labels (strings or integers), of shape `(n_samples)`:

```
>>> from sklearn import svm
>>> X = [[0, 0], [1, 1]]
>>> y = [0, 1]
>>> clf = svm.SVC()
>>> clf.fit(X, y)
SVC()
```

After being fitted, the model can then be used to predict new values:

```
>>> clf.predict([[2., 2.]])
array([1])
```

SVMs decision function (detailed in the [Mathematical formulation](#)) depends on some subset of the training data, called the support vectors. Some properties of these support vectors can be found in attributes `support_vectors_`, `support_` and `n_support_`:

```
>>> # get support vectors
>>> clf.support_vectors_
array([[0., 0.],
       [1., 1.]])
>>> # get indices of support vectors
>>> clf.support_
array([0, 1]...)
>>> # get number of support vectors for each class
>>> clf.n_support_
array([1, 1]...)
```

no seriously,
learn to read
docs



scikit-learn.org

...and general google-fu!

Cool, so... what's next?

coursera

Explore ▾

What do you want to learn?



Browse > Data Science > Machine Learning

Machine Learning Specialization

#BreakIntoAI with Machine Learning Specialization. Master fundamental AI concepts and develop practical machine learning skills in the beginner-friendly, 3-course program by AI visionary Andrew Ng

★★★★★ 4.9 9,029 ratings



Andrew Ng [+3 more instructors](#) **TOP INSTRUCTORS**

Enroll for Free
Starts Mar 19

Financial aid available

155,216 already enrolled



StatQuest with Josh Starmer ✓

@statquest 897K subscribers 245 videos

Statistics, Machine Learning and Data Science can sometimes :

towards data science

coursera

Explore ▾

What do you want to learn?



Browse > Data Science > Machine Learning

Deep Learning Specialization

Become a Machine Learning expert. Master the fundamentals of deep learning and break into AI. Recently updated with cutting-edge techniques!

★★★★★ 4.9 128,789 ratings



Andrew Ng [+2 more instructors](#) **TOP INSTRUCTORS**

Enroll for Free
Starts Mar 19

Financial aid available

747,307 already enrolled



(and have lots of fun!)

Ethics

- Interpretability of the model
- Security
- Privacy
- Trustworthiness
- Biases
- Fairness
- Socioeconomic consequences
- Unintended consequences
- Trolley problem

??????



Thank you.





AI Bridge

Samuel Ren, Henry M Gunn High School

Jiaming Situ, Homestead High School

Houjun Liu, The Nueva School

Xin Liu, Professor, Computer Science. UC Davis

With acknowledgments to

AI Institute for Food Systems

Saratoga Public Library